

# Introduction to Biological Computing: Masters in Computational Methods in Ecology & Evolution

(with modules for MSc/Mres EEC, EA, EChange, Trop Bio  
& NHM)

Imperial College London, 2015–16

Samraat Pawar ([s.pawar@imperial.ac.uk](mailto:s.pawar@imperial.ac.uk))  
*(with inputs from many others!)*

October 15, 2015



# Contents

<b>0 Introduction</b>	<b>1</b>
0.1 Introduction . . . . .	1
0.2 The goals of this course . . . . .	2
0.2.1 Assessment of practicals . . . . .	3
0.3 Conventions used in this document . . . . .	4
0.4 To IDE or not to IDE? . . . . .	4
<b>1 Introduction to UNIX and Linux</b>	<b>7</b>
1.1 What is UNIX? . . . . .	7
1.2 Why UNIX? . . . . .	7
1.3 UNIX directory structure . . . . .	8
1.4 Meet the UNIX shell . . . . .	8
1.5 sudo, and installing/removing software . . . . .	9
1.6 Basic UNIX commands . . . . .	10
1.7 Building your coursework directory structure . . . . .	10
1.8 Command arguments . . . . .	11
1.9 Redirection and pipes . . . . .	12
1.10 Wildcards . . . . .	12
1.11 Using grep . . . . .	13
1.12 Finding files with find . . . . .	15
1.13 Tweaking your UNIX and bash . . . . .	16
1.13.1 Opening nautilus from terminal . . . . .	16
1.13.2 Enabling or improving autocomplete in terminal . . . . .	17
1.14 Practical: Make sure the basics work . . . . .	17
1.15 Readings & Resources . . . . .	18
<b>2 Version control with Git</b>	<b>19</b>
2.1 What is Version Control? . . . . .	19
2.2 Why Version Control? . . . . .	20
2.3 git . . . . .	20
2.4 Your first repository . . . . .	20
2.5 git commands . . . . .	21
2.5.1 git command structure . . . . .	21
2.6 Ignoring Files . . . . .	22
2.7 Removing files . . . . .	23
2.8 Accessing history of the repository . . . . .	23
2.9 Reverting to a previous version . . . . .	24
2.10 Branching . . . . .	24
2.11 Running git commands on a different directory . . . . .	25

2.12	Running git commands on multiple repositories at once . . . . .	25
2.13	Wrapping up . . . . .	26
2.14	Practical . . . . .	26
2.15	Readings & Resources . . . . .	26
<b>3</b>	<b>Advanced UNIX: Shell scripting</b>	<b>27</b>
3.1	Shell scripting: What and Why . . . . .	27
3.2	Scripting: How . . . . .	27
3.3	Your first shell script . . . . .	28
3.4	A useful shell-scripting example . . . . .	28
3.5	Variables in shell scripting . . . . .	29
3.6	Some more Examples . . . . .	30
3.7	Practical . . . . .	30
3.8	Readings & Resources . . . . .	31
<b>4</b>	<b>Using L<sup>A</sup>T<sub>E</sub>X for scientific documents</b>	<b>33</b>
4.1	What's L <sup>A</sup> T <sub>E</sub> X? . . . . .	33
4.2	Why L <sup>A</sup> T <sub>E</sub> X? . . . . .	33
4.2.1	Limitations of L <sup>A</sup> T <sub>E</sub> X . . . . .	34
4.3	Installing L <sup>A</sup> T <sub>E</sub> X . . . . .	34
4.4	A first L <sup>A</sup> T <sub>E</sub> X example . . . . .	34
4.4.1	A bash script to compile L <sup>A</sup> T <sub>E</sub> X . . . . .	36
4.5	A brief L <sup>A</sup> T <sub>E</sub> X tour . . . . .	36
4.6	L <sup>A</sup> T <sub>E</sub> X templates . . . . .	37
4.7	Typesetting math . . . . .	38
4.8	A few more tips . . . . .	38
4.9	Readings & Resources . . . . .	39
4.10	Practical . . . . .	39
<b>5</b>	<b>Biological Computing in Python I</b>	<b>41</b>
5.1	Chapter 1's UNIX question? . . . . .	41
5.2	Outline of the python module . . . . .	41
5.3	Why python? . . . . .	41
5.3.1	The Zen of python . . . . .	42
5.4	Installing python . . . . .	42
5.5	python warmup . . . . .	43
5.5.1	Variable types . . . . .	43
5.5.2	python operators . . . . .	44
5.5.3	Assigning and manipulating variables . . . . .	44
5.6	python data types and data structures . . . . .	45
5.6.1	Lists . . . . .	45
5.6.2	Tuples . . . . .	46
5.6.3	Sets . . . . .	46
5.6.4	Dictionaries . . . . .	47
5.6.5	Copying mutable objects . . . . .	48
5.6.6	python with strings . . . . .	49
5.7	python Input/Output . . . . .	49
5.8	Writing python code . . . . .	51
5.8.1	Writing python functions (or modules) . . . . .	52
5.8.2	Components of the python function . . . . .	52

5.8.3	Variable scope . . . . .	55
5.8.4	Running python code from terminal or bash shell . . . . .	56
5.9	Control statements . . . . .	56
5.9.1	Control flow exercises . . . . .	57
5.10	Loops . . . . .	59
5.10.1	List comprehensions . . . . .	59
5.11	Practical . . . . .	60
5.12	Readings and Resources . . . . .	61
<b>6</b>	<b>Biological Computing in Python II</b>	<b>63</b>
6.1	ipython . . . . .	63
6.2	Functions, Modules, and code compartmentalization . . . . .	65
6.2.1	Importing Modules . . . . .	65
6.3	Python packages . . . . .	65
6.4	Regular expressions in python . . . . .	66
6.4.1	Some RegExercises . . . . .	70
6.4.2	Important re functions . . . . .	71
6.5	Numerical computing in python . . . . .	71
6.5.1	A numerical integration example . . . . .	76
6.6	Practical . . . . .	76
6.6.1	Align DNA sequences . . . . .	76
6.6.2	Blackbirds problem . . . . .	76
6.7	Readings and Resources . . . . .	76
<b>7</b>	<b>Additional Python topics, and a Long practical</b>	<b>79</b>
7.1	Using python to build workflows . . . . .	79
7.2	Using subprocess . . . . .	79
7.2.1	Running R . . . . .	80
7.3	Long Practical – Building a workflow . . . . .	81
7.3.1	Background and Objectives . . . . .	81
7.3.2	The Data . . . . .	81
7.3.3	The Models . . . . .	81
7.3.4	Fitting data to the models . . . . .	82
7.3.5	The Workflow . . . . .	82
7.3.6	Patching together your workflow components . . . . .	84
7.3.7	Submission . . . . .	84
7.4	Errors in your python code . . . . .	84
7.4.1	Unit testing . . . . .	85
7.4.2	Debugging . . . . .	87
7.4.3	Debugging with breakpoints . . . . .	89
7.5	Profiling in Python . . . . .	89
7.6	Practical . . . . .	91
7.6.1	Lotka-Volterra model problem . . . . .	91
7.6.2	Missing oaks problem! . . . . .	92
7.7	Databases and python . . . . .	92
7.7.1	python with SQLite . . . . .	94
7.8	Practical . . . . .	100
7.8.1	Using os problem 1 . . . . .	101
7.8.2	Using os problem 2 . . . . .	101
7.9	Readings and Resources . . . . .	101

<b>8 Introduction to R</b>	<b>103</b>
8.1 Outline of the the R module . . . . .	103
8.2 What is R? . . . . .	103
8.3 Would you ever need anything other than R? . . . . .	103
8.4 Installing R . . . . .	104
8.5 Getting started . . . . .	104
8.6 Useful R commands . . . . .	106
8.7 Some Basics . . . . .	106
8.7.1 Variable names and Tabbing . . . . .	106
8.7.2 R likes E Notation . . . . .	107
8.7.3 Operators . . . . .	107
8.7.4 When things go wrong . . . . .	108
8.7.5 Types of parentheses . . . . .	108
8.8 Data types . . . . .	108
8.8.1 Vectors . . . . .	108
8.8.2 Matrices and arrays . . . . .	109
8.8.3 Data frames . . . . .	109
8.8.4 Lists . . . . .	111
8.9 Variable Types, Type Conversion and Special Values . . . . .	111
8.10 Creating and Manipulating Data structures . . . . .	112
8.10.1 Sequences . . . . .	112
8.10.2 Strings and Pasting . . . . .	112
8.10.3 Indices and Indexing . . . . .	113
8.10.4 Recycling . . . . .	114
8.10.5 Basic vector-matrix operations . . . . .	114
8.11 Your R Analysis Workflow . . . . .	114
8.11.1 The Working Directory . . . . .	114
8.12 The R analysis workflow . . . . .	115
8.13 Importing and Exporting Data . . . . .	116
8.14 Writing Functions . . . . .	117
8.15 Control statements . . . . .	117
8.16 Running R code . . . . .	118
8.17 Useful R Functions . . . . .	119
8.17.1 Mathematical . . . . .	119
8.17.2 Strings . . . . .	119
8.17.3 Statistical . . . . .	119
8.17.4 Random number distributions . . . . .	119
8.18 Packages . . . . .	120
8.19 Practical . . . . .	120
8.20 Readings . . . . .	120
<b>9 Plotting and graphics in R</b>	<b>123</b>
9.1 Basic plotting . . . . .	123
9.1.1 Scatter Plot . . . . .	124
9.1.2 Histograms . . . . .	125
9.1.3 Subplots . . . . .	125
9.1.4 Overlaying plots . . . . .	125
9.1.5 Exercise . . . . .	126
9.1.6 Saving your graphics . . . . .	126

9.1.7	Boxplots . . . . .	126
9.1.8	Practical 2.1 . . . . .	126
9.2	Publication-quality figures in R . . . . .	127
9.2.1	Basic graphs with <code>qplot</code> . . . . .	129
9.2.2	Various <code>geom</code> . . . . .	132
9.2.3	Practical 2.2 . . . . .	133
9.2.4	Advanced plotting: <code>ggplot</code> . . . . .	133
9.2.5	Case study 1: plotting a matrix . . . . .	135
9.2.6	Case study 2: plotting two dataframes . . . . .	137
9.2.7	Case study 3: annotating the plot . . . . .	138
9.2.8	Case study 4: mathematical display . . . . .	139
9.3	Readings . . . . .	139
<b>10</b>	<b>Advanced topics in R</b> . . . . .	<b>143</b>
10.1	Vectorization . . . . .	143
10.2	Breaking out of loops . . . . .	144
10.3	Using <code>next</code> . . . . .	145
10.4	“Catching” errors . . . . .	145
10.5	Generating Random Numbers . . . . .	146
10.6	Debugging . . . . .	146
10.7	Practical . . . . .	147
10.8	Launching/Running R in batch mode . . . . .	148
10.9	Accessing databases using R . . . . .	148
10.10	Building your own R packages . . . . .	148
10.11	<code>Sweave</code> . . . . .	148
10.12	Practical (CMEE Extra Credit) . . . . .	149
10.13	Practical (CMEE Extra Credit) . . . . .	149
10.14	R Module Wrap up . . . . .	150
10.14.1	Some comments and suggestions . . . . .	150
10.14.2	CMEE weekly coursework submission . . . . .	150
10.15	Readings . . . . .	150



# Chapter 0

## Introduction

It is hard for me to say confidently that, after fifty more years of explosive growth of computer science, there will still be a lot of fascinating unsolved problems at peoples' fingertips, that it won't be pretty much working on refinements of well-explored things. Maybe all of the simple stuff and the really great stuff has been discovered. It may not be true, but I can't predict an unending growth. I can't be as confident about computer science as I can about biology. Biology easily has 500 years of exciting problems to work on, it's at that level.

---

— Donald Knuth

### 0.1 Introduction

This document contains the content of the modules on Biological Computing in the MSc/MREs Courses on Computational Methods in Ecology and Evolution (CMEE) at Silwood Park, Department of Life Sciences, Imperial College London. These notes are also available to other Masters courses, including EA, EEC, EEChange, Trop Bio and NHM, and the chapters on R will be necessary for these courses. Please look up your respective class handbooks to determine when the computing week(s) are scheduled in your particular course (e.g., the R module, sub-modules of which are taught across all the courses).

Each chapter corresponds to a module, or part of a module, and will typically be covered within a single week of the first nine weeks of the Autumn semester.

This document is accompanied by data and code on which you can practice your skills in your own time and during the practical sessions. These materials are available (and will be updated regularly) at: <https://bitbucket.org/mhasoba/cmee2015masterrepo>.

It is important that you work through the exercises and problems in each chapter/section. This document does not tell you every single thing you need to know to perform the exercises in it. In

programming and computing, you learn faster by trying to solve problems on your own (including computer crashes!), often by liberally googling them!

You will be provided guidelines for what makes good or efficient solutions to the computing exercises. Later, when you have submitted your exercises and practicals (only relevant to the CMEE students, but the rest should have a go!), you will be provided solutions.

Also, every time a mysterious, geeky-sounding term like “relative path” or “version control” appears, please google it!

## 0.2 The goals of this course

The goal of this course is to teach you to become (or at least show you the path towards becoming) a competent quantitative biologist. A large part of this involves learning computer programming. Why do biologists need to write write computer programs? Here are my top reasons:

- Short of fieldwork, programs can do anything (that can be specified). In fact, even fieldwork, if you could program a robot to do it for you!
- As such, no software is typically available to perform exactly the analysis you are planning. You should be concerned if you are trying to shoehorn your data into methods that don’t quite seem right.
- Programming permits success despite complexity, through modularization and precise specification. Biological problems and datasets are some of the most complicated imaginable.
- Modularity – programming allows you to break up your complex analysis in smaller pieces, yet keep all the pieces in a single, functional analysis.
- Reproducibility – you (or someone else) can just re-run the code to reproduce your analysis. This is also the key to maintaining scientific accountability, integrity, and accuracy.
- Organised thinking – writing code requires you to do this!
- Career prospects – good, scientific coders are in short supply in all fields, but most definitely in biology!

There are several hundred programming languages currently available – which ones should a biologist choose? Ideally, in scientific computing in biology, one would like to know:

1. A fast, compiled (or semi-compiled) “procedural” language (e.g., C) – not so important for most of you, and beyond the scope of this course
2. A modern, easy-to-write, interpreted (or semi-compiled) language that is still quite fast, like python
3. A mathematical/statistical software with programming and graphing capabilities like R (also, mathematica, MATLAB)

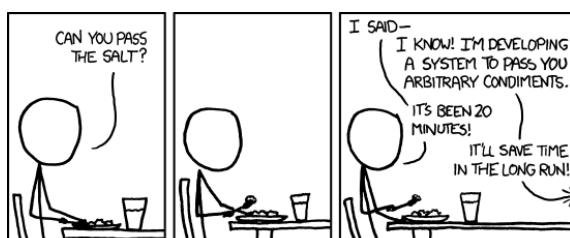


Figure 1: But not every task needs to be converted to a computer program — you will learn to decide when to and when not to write a computer program! <http://xkcd.com/974/>

All these because one language doesn't fit all purposes! Therefore you will learn a few different languages in this course — hopefully, just the right number!

Our goal is to teach you not just programming, but also good computing practices. In this course, you will write plenty of code, deal with different data files, and produce text and graphic outputs. You will learn to keep your project and coursework organized in logical, efficient, error-free and reproducible *workflows* (that's a mouthful, but an important mouthful).

### 0.2.1 Assessment of practicals

If you are a CMEE student, we will assess both, your practical computing work itself, and whether you are following good programming and workflow practices on a weekly basis. This will be done using scripts — yes, that's how hard-core we are — we will assess your scripts using scripts! I have written a python script that checks whether your weekly (and version-controlled) directories are neat and organized in a logical workflow, and whether all the scripts run correctly with the expected inputs and outputs, starting from week 1 (Chapter 1). Specifically, as an example towards learning good workflow practices, you will keep all your coursework code, data inputs and results outputs organized in separate directories named `Code`, `Data`, `Results` (or equivalent) respectively.

More details on this will be given when the weekly modules actually start. The assessment script will then record a log file that summarizes all the issues found in your workflows, which will be emailed to you by the middle of the week subsequent to the week you submitted your weekly practical work (usually on a wednesday). A manually-written assessment of your overall performance will be sent at the end of the 9 weeks.

Note that practicals in the weeks/modules not in these notes (e.g., GIS, Genomics) may also be included in the assessment (you will be given clear instructions if that is the case).

Here is the marking scheme I will use:

1. You would get 100% if,
  - (a) All the in-class scripts were in place (in the code directory in the respective week's directory) and functional when run on my computer
  - (b) All the assigned practicals / problems are complete and functional, and give the right answers
  - (c) The scripts are all up to the mark in terms of internal documentation and commenting
  - (d) There is a neat `readme` file detailing all the scripts and what they do week 1 onwards.



EVERY CHANGE BREAKS SOMEONE'S WORKFLOW.

Figure 2: Logical workflows are important, but don't get addicted to yours!  
<http://xkcd.com/1172/>

2. For every script that gives a syntax error, 5 pts deducted, and for every script that gives an error because of wrong path (e.g., absolute) assignment, 2 pts deducted.
3. For every missing assigned practical/problem, 10 pts deducted
4. For every assigned practical/problem, 5 pts deducted for wrong answer if applicable (that is, script runs without error, but gives wrong numerical or text output).
5. For every missing `readme` file, 1 pt deducted
6. For every extra-credit question completed, 2.5 pts added.
7. Apart from these objective criteria, I will exercise my judgement to deduct more marks if the weekly directory structure is disorganized, the code inadequately commented or insufficiently documented, or something else. You will get feedback if these issues needed to be addressed in the final written assessment.

In the above marking scheme, “in-class scripts” are the ones that were given to you to practice with, which you only had to reproduce without error; the “assigned practicals / problems” are the assignments/problems you have been given that involve the writing of new scripts or the modification of existing ones.

### 0.3 Conventions used in this document

Throughout this document, directory paths will be specified in UNIX (Mac, Linux) style, using `/` instead of the `\` used in Windows. Also, in general, we will be using *relative paths* throughout the exercises and practicals (more on this later, but google it!).

You will find all command line/console arguments, code snippets and output in coloured boxes like this:

```
$ ls
```

The specific prompt (`$` in this case, belonging to the UNIX terminal) will vary with the programming language/console (`$` for UNIX, `>>>` for Python, `>` for R, etc.)

You will type the commands/code that you see in such boxes into the relevant command line (or copy-paste, but not recommended!). I have aimed to make the content of this module computer platform (Mac, Linux or PC) independent. Also note that:

- \* Lines marked with a star like this will be specific instructions for you to follow

### 0.4 To IDE or not to IDE?

As you embark on your journey to becoming a competent practitioner of biological computing, you will be faced with a hamletian question: “To IDE or not to IDE”. Ok, maybe not that dramatic or hamletian... Anyway, an interactive Development Environment (IDE) is a text editor with frills that can make life easy by auto-formatting code, running code, allowing a graphic view of the workspace (your active functions, variables, etc.), graphic debugging and profiling (you will see

these delightful things later), and allowing integrated version control (e.g., using `git`). You will benefit a lot if you use a good code editor that can also offer an IDE (e.g., emacs, vim, geany). At the very least, your IDE should offer:

- Autoindentation
- Automatic code wrapping (e.g., keeping lines <80 characters long)
- Syntax highlighting (e.g., commands vs. variables)
- Code folding (fold large blocks of code, say an entire function or loop)
- Keyboard control of commenting/uncommenting, code wrapping, etc.
- Embedded terminal / shell / commandline console
- Sending commands to terminal / shell

And if you end up using multiple programming languages, you will want an IDE that can handle them. For example, RStudio cannot handle more than a fixed set of 2-3 languages. I use geany, which has many plugins that make multi-language (multilingual?) code development much easier. I would also recommend vim or emacs, which have a steeper learning curve, but are very powerful once you have mastered them.

Alright, full steam ahead then!



# Chapter 1

## Introduction to UNIX and Linux

### 1.1 What is UNIX?

UNIX is an operating system developed in the 1970s by AT&T programmers (notably Brian Kernighan and Dennis Ritchie, fathers of C). It is multi-user and network-oriented by design. It used plain text files for storing data and a strictly hierarchical file system.

It is a machine-independent OS. Linux and Mac OS are Unix-like (or UN\*X or \*nix) operating systems. Ubuntu is a Linux distribution.

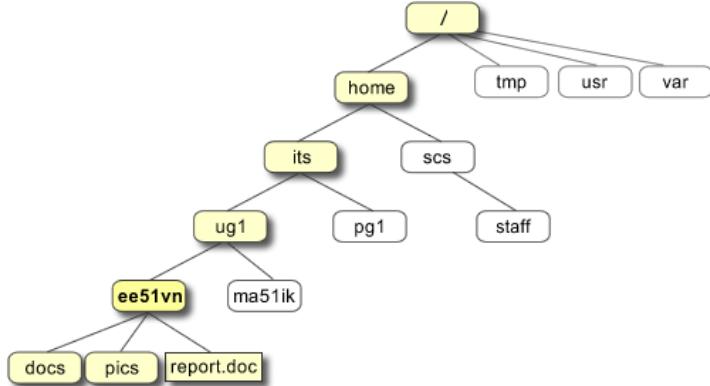
### 1.2 Why UNIX?

Here are some very good reasons:

- It was designed for developing code and storing data
- Scores of small programs available to perform simple tasks – can be combined easily
- Easy to automate tasks (e.g., using shell scripts)
- Multi-user (multiple users can log in concurrently use computer)
- Multi-tasking (can perform many tasks at the same time)
- Network-ready (easy to communicate between computers)
- Robust, stable, secure (very few UNIX viruses and malware — I have never encountered one!)
- Free and open source!

See <http://www.whylinuxisbetter.net/> to aid your brain-washing (cleaning?)!

### 1.3 UNIX directory structure



<http://pathanruet.files.wordpress.com/2012/05/unix-tree.png>

The UNIX directory structure is hierarchical, with a single tree starting from the “root” /. This is quite unlike Windows or MS-DOS, where there are separate trees for disk partitions, removable media, network, etc. The key UNIX directories are:

/	Is the “root” directory
/bin	Contains basic programs
/etc	Contains configuration files
/dev	Contains files connecting to devices (keyboard, mouse, screen, etc.)
/home	Your home directory – this is where you will usually work
/tmp	Contains Temporary files

### 1.4 Meet the UNIX shell

The shell (or terminal) is a text command processor to interface with the Operating System’s “kernel”. We will use the popular (yes, it’s popular!) bash shell.

- \* To launch bash shell, do **Ctrl + Alt + t** (or use **Meta** key)

OK, so you have met the shell. Note that:

- The shell automatically starts in your home directory `/home/yourname/`, also called `~` (important to remember!)
- Use the Tab key – very handy (try `ls` with Tab Tab)
- Navigate commands you typed using the up/down arrows

Other useful keyboard shortcuts are:

Ctrl + A	Go to the beginning of the line
Ctrl + E	Go to the end of the line
Ctrl + L	Clear the screen
Ctrl + U	Clear the line before cursor position
Ctrl + K	Clear the line after the cursor
Ctrl + C	Kill whatever you are running
Ctrl + D	Exit the current shell
Ctrl + right arrow	Move cursor forward one word
Ctrl + left arrow	Move cursor backward one word

## 1.5 sudo, and installing/removing software

You can install software in your `/home` directory. In UNIX you originally had to login as `root` (administrator). But in Ubuntu, it is sufficient to add `sudo` (super user `do`) in front of a command:

```
sudo apt-get install geany geany-plugins geany-plugin-latex  
geany-plugin-addons
```

You can install anything that is in the Ubuntu “repository”. Let’s try installing something else — a weather indicator that I think really works very well.

But here you have to first install the repository:

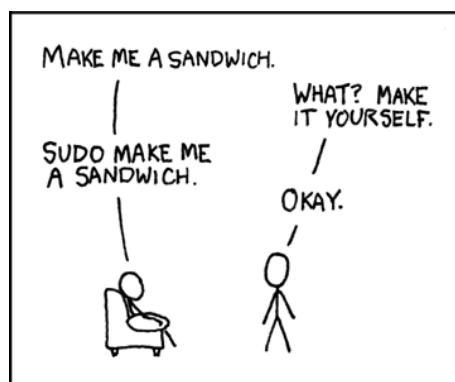
```
$ sudo add-apt-repository ppa:atareao/atareao
```

This will ask you to first approve the addition of this repository. Then update the repository packages,

```
$ sudo apt-get update
```

And then install,

```
$ sudo apt-get install my-weather-indicator
```



<http://xkcd.com/149/>

You can also easily remove software by, well, using the `remove` command!

```
$ sudo apt-get remove indicator-messages
```

This will get rid of the evolution mail indicator — very unlikely that you will use evolution!

## 1.6 Basic UNIX commands

<code>man [COMMAND]</code>	Show help page of a command.
<code>whoami</code>	Display your user-name.
<code>pwd</code>	Show the current directory.
<code>ls</code>	List the files in the directory.
<code>cd [DIRNAME]</code>	Change directory.
<code>cd ..</code>	Move one directory up.
<code>cd /</code>	Go to the root directory.
<code>cd ~ or just cd</code>	Go to the home directory.
<code>cp [FROM] [TO]</code>	Copy a file or directory non-recursively (what's this?).
<code>mv [FROM] [TO]</code>	Move or rename a file or directory.
<code>touch [FILENAME]</code>	Create an empty file.
<code>echo "My string"</code>	Print a string (here, "My string").
<code>rm [TOREMOVE]</code>	Remove a file or directory non-recursively.
<code>mkdir [DIRNAME]</code>	Create a directory.
<code>rmdir [DIRNAME]</code>	Remove an empty directory.
<code>wc [FILENAME]</code>	Count the number of lines and words in a file.
<code>sort [FILENAME]</code>	Sort the lines of a file and print result.
<code>uniq</code>	Shows only unique elements of a list.
<code>cat [FILENAME]</code>	Print the file on the screen.
<code>less [FILENAME]</code>	Progressively print a file on the screen ("q" to exit).
<code>head [FILENAME]</code>	Print the first few lines of a file.
<code>tail [FILENAME]</code>	Print the last few lines of a file.
<code>history</code>	Show the last commands you typed.
<code>date</code>	Print current date.
<code>file [FILENAME]</code>	Determine the type of a file.
<code>passwd</code>	Change user password.
<code>chmod [FILENAME]</code>	Change file permissions.

## 1.7 Building your coursework directory structure

It is time to start building your CMEE coursework directory structure. Please follow these rules:

- Do all your work in CMEECourseWork, located in a suitable place in your /home (make mama proud, keep your home organized!)
  - Each week's coursework should be in its respective directory: CMEECourseWork/Week1, CMEECourseWork/Week2, etc
  - Each week's directory will contain directories called Code, Data, etc (later)
  - You will bring CMEECourseWork and all it's contents under version control using Git (later)
- \* OK, `mkdir` your CMEECourseWork directory now.

Starting by creating CMEECourseWork). *Don't forget to use the tab key freely!*. Then,

```
$ mkdir Week1
$ cd Week1
$ mkdir sandbox
$ cd Sandbox
bash: cd: Sandbox: No such file or directory
$ cd ..
$ rm sandbox
rm: cannot remove `sandbox/`: Is a directory
$ mv sandbox Sandbox # OR, "rm -r sandbox" (careful with the -r option!)
```

Note the hash mark # above — anything after a # is ignored (so you can use it for commenting).

```
$ cd Sandbox
$ pwd
$ ls
$ touch TestFile # OR, "touch TestFile.txt"
$ ls
$ mv TestFile TestFile2
$ rm TestFile2
```

You could have made your project directories and subdirectories in one swoop by using the -p option of mkdir:

```
$ mkdir -p CMEECourseWork/Week1/{Data,Code,Sandbox}
```

## 1.8 Command arguments

Most UNIX commands accept arguments that modify their behavior. E.g., `ls -l` (`ls` “minus”l) lists the files in longer format. Some useful arguments:

<code>cp -r [DIR1] [DIR2]</code>	Copy a directory recursively (i.e., including all the sub-directories and files).
<code>rm -i [FILENAME]</code>	Remove a file, but asks first (for safety).
<code>rm -r [DIR]</code>	Remove a directory recursively (i.e., including all the sub-directories and files).
<code>ls -a</code>	List all files, including hidden ones.
<code>ls -h</code>	List all files, with human-readable sizes (Mb, Gb).
<code>ls -l</code>	List all files, long format.
<code>ls -S</code>	List all files, order by size.
<code>ls -t</code>	List all files, order by modification time.
<code>ls -1</code>	List all files, one file per line.
<code>mkdir -p Dir1/Dir2/Dir3</code>	Create the directory Dir3 and Dir1 and Dir2 if they do not already exist.
<code>sort -n</code>	Sort all the lines, but use numeric values instead of dictionary (i.e., 11 follows 2).

## 1.9 Redirection and pipes

Output of programs can also be “redirected” to a file:

- >      Redirect output from a command to a file on disk. If the file already exists, it will be overwritten.
- >>     Append the output from a command to a file on disk. If the file does not exist, it will be created.

Examples (make sure you are in Week1/Sandbox):

```
$ echo "My first line." > test.txt
$ cat test.txt
$ echo "My second line" >> test.txt
$ ls / >> ListRootDir.txt
$ cat ListRootDir.txt #That's cool! Why?
```

We can also concatenate commands using “pipes” with “|” e.g., to count how many files are in root (/) directory:

```
$ ls / | wc -l #look up "info wc"
```

## 1.10 Wildcards

We can use wildcards to find files based on their names (again, in Week1/Sandbox !):

```
$ mkdir TestWild
$ cd TestWild
$ touch File1txt
$ touch File2.txt
$ touch File3.txt
$ touch File4.txt
$ touch File1.csv
$ touch File2.csv
$ touch File3.csv
$ touch File4.csv
$ touch Anotherfile.csv
$ touch Anotherfile.txt
$ ls
$ ls | wc -l
```

We will use the following wildcards:

- |       |   |
|-------|---|
| ?     | Any single character, except a leading dot (hidden files).    |
| *     | Zero or more characters, except a leading dot (hidden files). |
| [A-Z] | Define a class of characters (e.g., upper-case letters).      |

Now let's try to find the files using wildcards:

```
$ ls *
$ ls File*
```

```
$ ls *.txt
$ ls File?.txt
$ ls File[1-2].txt
$ ls File[!3].*
```

## 1.11 Using grep

grep is a command that matches strings in a file (why is this useful?). It is based on regular expressions (more on this later). Let's explore some basic usage of grep. For a test file let's download a list of protected species from the UN website (to Sandbox):

```
$ wget http://www.cep.unep.org/pubs/legislation/spawannxs.txt #Cool!
$ head -n 50 spawannxs.txt #You will see "head" in R as well
```

Now,

```
$ mkdir ../Data #Note the relative path "../"
$ mv spawannxs.txt ../Data/
$ cd ../Data
$ head -n 50 spawannxs.txt
```

Note that now you have a Data directory.

OK, what about falcons?

```
$ grep Falco spawannxs.txt
Falconidae Falco femoralis septentrionalis
Falconidae Falco peregrinus
Falconidae Polyborus plancus
Falconidae Falco columbarius
```

Using `-i` make the matching case-insensitive:

```
$ grep -i Falco spawannxs.txt
Order: FALCONIFORMES
Falconidae Falco femoralis septentrionalis
Falconidae Falco peregrinus
Falconidae Polyborus plancus
Order: FALCONIFORMES
Order: FALCONIFORMES
Order: FALCONIFORMES
Falconidae Falco columbarius
```

Now let's find the beautiful “Ara” macaws:



But this poses a problem (what is the problem?):

```
$ grep -i ara spawannxs.txt
Flacourtiaceae Banaras      vanderbiltii
Order: CHARADRIIFORMES
Charadriidae Charadrius melodus
Psittacidae Amazona     arausica
Psittacidae Ara       macao
Dasyprotidae Dasyprocta guamara
Palmae      Syagrus (= Rhyticoccos) amara
Psittacidae Ara       ararauna
Psittacidae Ara       chloroptera
Psittacidae Arao      manilata
Mustelidae Eira       barbara
Order: CHARADRIIFORMES
```

We can solve this by specifying `-w` to match only full words:

```
$ grep -i -w ara spawannxs.txt
Psittacidae Ara      macao
Psittacidae Ara     ararauna
Psittacidae Ara     chloroptera
```

And also show line(s) after the one that was matched, we can use `-A x`, where `x` is number of lines to use:

```
$ grep -i -w -A 1 ara spawannxs.txt
Psittacidae Ara      macao

--
Psittacidae Ara     ararauna
Psittacidae Ara     chloroptera
Psittacidae Arao    manilata
```

Similarly, `-B` shows the lines before:

```
$ grep -i -w -B 1 ara spawannxs.txt
Psittacidae Amazona      vittata
Psittacidae Ara      macao
--
Psittacidae Amazona      ochrocephala
Psittacidae Ara      ararauna
Psittacidae Ara      chloroptera
```

Use `-n` to show the line number of the match:

```
$ grep -i -w -n ara spawannxs.txt
216:Psittacidae Ara      macao
461:Psittacidae Ara      ararauna
462:Psittacidae Ara      chloroptera
```

To print all the lines that do not match a pattern, use `-v`:

```
$ grep -i -w -v ara spawannxs.txt
```

To match one of several strings, use `grep "string1|string2" file`. `grep` can be used on multiple files, all files, using wildcards for filenames, etc – explore as and when you need.

## 1.12 Finding files with find

It's easy to find files in UNIX using `find`! Let's test it (make sure you are in `Sandbox`, not `Data!`)

```
$ mkdir TestFind
$ cd TestFind
$ mkdir -p Dir1/Dir11/Dir111 #what does -p do?
$ mkdir Dir2
$ mkdir Dir3
$ touch Dir1/File1.txt
$ touch Dir1/File1.csv
$ touch Dir1/File1.tex
$ touch Dir2/File2.txt
$ touch Dir2/file2.csv
$ touch Dir2/File2.tex
$ touch Dir1/Dir11/Dir111/File111.txt
$ touch Dir3/File3.txt
```

Now find particular files:

```
$ find . -name "File1.txt"
./Dir1/File1.txt
```

Using `-iname` ignores case, and you can use wildcards:

```
$ find . -iname "fi*.csv"
./Dir1/File1.txt
./Dir1/Dir11/Dir111/File111.txt
./Dir3/File3.txt
./Dir2/File2.txt
```

You can limit the search to exclude sub-directories:

```
$ find . -maxdepth 2 -name "*txt"
./Dir1/File1.txt
./Dir3/File3.txt
./Dir2/File2.txt
```

You can exclude certain files:

```
$ find . -maxdepth 2 -not -name "*txt"
.
./Dir1
./Dir1/File1.tex
./Dir1/File1.csv
./Dir1/Dir11
./Dir3
./Dir2
./Dir2/File2.tex
./Dir2/File2.csv
```

To find only directories:

```
$ find . -type d
.
./Dir1
./Dir1/Dir11
./Dir1/Dir11/Dir111
./Dir3
./Dir2
```

## 1.13 Tweaking your UNIX and bash

There are numerous ways in which you can tweak your Linux/UNIX environment and bash/terminal to your likes. For example, see <http://www.howtogeek.com/tag/ubuntu/ubuntu-tips/>.

But be careful, it can be addictive, and sometime dangerous to your system's stability!. For example, here are a couple of tweaks that I really find useful:

### 1.13.1 Opening nautilus from terminal

In terminal you can enter “f” to open nautilus in current directory by doing the following. Open your `.bashrc` for editing:

```
$ sudo gedit ~/.bashrc
```

Then add to the last line (type it, don't copy and paste!):

```
alias f='nautilus .'
```

Then restart terminal or in current terminal:

```
$ source ~/.bashrc
```

### 1.13.2 Enabling or improving autocomplete in terminal

What happens when you use up and down keys in terminal? If nothing, then you need to enable reverse searching history. To do so, open /etc/inputrc

```
$ sudo open /etc/inputrc
```

Then, add the following to it:

```
"\eOA": history-search-backward\\
"\e[A": history-search-backward\\
"\eOB": history-search-forward\\
"\e[B": history-search-forward\\
"\eOC": forward-char\\
"\e[C": forward-char\\
"\eOD": backward-char\\
"\e[D": backward-char\\
```

]

If this doesn't work, google it! Also, see the solution at <http://askubuntu.com/questions/59846/bash-history-search-partial-up-arrow>

Then close current terminal, open new one, and try up and down keys again.

## 1.14 Practical: Make sure the basics work

### 1. Some instructions:

Review (especially if you got lost along the way) and make sure you can run and understand all the commands and get the expected outputs we have covered today.

Make sure you have your directory organized with Data and Sandbox with the necessary files, under CMEECourseWork/Week1.

Along with the completeness of the practicals/exercises themselves, you will be marked on the basis of how complete and well-organized your directory structure and content is – in all coming weeks as well.

2. Here is a more complicated bash command using two pipes (you are not expected to include the answer to this one as part of your weekly submission):

```
$ find . -type f -exec ls -s {} \; | sort -n | head -10
```

What does this command do (Hint: try it on the test directories and files we created in Sandbox)?

Note that along with the `man` command, you can use the internet to get help on practically everything about UNIX!

3. In the directory `UNIX/Data/fasta` you find some FASTA files. These files have an header starting with `>` followed by the name of the sequence and other metadata. Starting from the second line, we have the sequence data. Write a file called `UnixPrac1.txt` with UNIX shell commands that do the following (number each command with a hashed comment like so – # 1, # 2, etc):

- (a) Count how many lines are in each file
- (b) Print everything starting from the second line for the *E. coli* genome
- (c) Count the sequence length of this genome
- (d) Count the matches of a particular sequence, “ATGC” in the genome of *E. coli* (hint: Start by removing the first line and removing newline characters)
- (e) Compute the AT/GC ratio

Save `UnixPrac1.txt` in the `Code` directory. Please make sure that each command calls the data from the `Data` directory! Do not write any of the above as shell scripts (that's not been covered yet; see Chapter 3) — each one should be a single line solution made of (potentially piped together) UNIX commands.

**It will if you have (judicious) comments in any of your script files. But I won't penalize you if you haven't put in comments in the first week in practicals where I don't explicitly ask for them. From the first Python week (Chapter 5) onwards, I will penalize you if you don't properly document and comment code (more on this next week), even if I haven't explicitly asked you**

## 1.15 Readings & Resources

- Lots of UNIX tutorials out there. Try <http://software-carpentry.org/lessons.html> (Chapter “shell”). ( watch video tutorials or read pdfs)
- IC library gives you with access to several e-books on UNIX, some specific to Ubuntu. Browse or search and find a good intro book: [http://imp-primo.hosted.exlibrisgroup.com/primo\\_library/libweb/action/search.do](http://imp-primo.hosted.exlibrisgroup.com/primo_library/libweb/action/search.do)
- There are also paper books – again, search library website
- List of UNIX commands along with man page: [www.oreillynet.com/linux/cmd/](http://www.oreillynet.com/linux/cmd/)

# Chapter 2

## Version control with Git

### 2.1 What is Version Control?

Version control, also known as revision control or source control, is the management of changes to documents, computer programs, large web sites, and other collections of information in an automated way.

Any project (collections of files in directories) under version control has changes and additions/deletions to its files and directories recorded and archived over time so that you can recall specific versions later. With version control of biological computing projects, you can:

- record of all changes made to a set of files and directories, including text (usually ascii) data files, so that you can access any previous version of the files
- branch (and merge) new projects
- “roll back” data, code, documents (but not formats like \*.docx though!)

Version control is embedded in various word processors and spreadsheets, e.g., Google Docs and Sheets

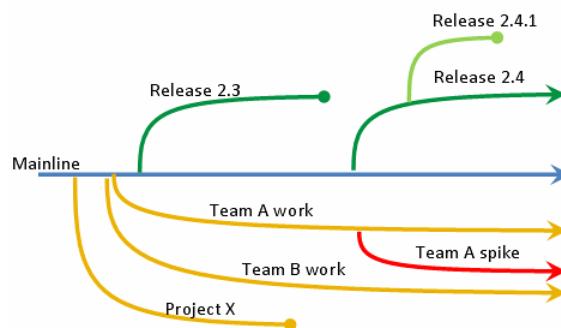


Figure 2.1: A general idea of how version control works.

## 2.2 Why Version Control?

Frist day after the project is assigned ..

```
mak@company $ mkdir proj
mak@company $ ls proj
index.html
```

After a week .. !!!!!

```
mak@company $ ls proj
header.php          header1.php    header2.php
header_current.php index.html      index.html.bkp
index.html.old
```

After a fortnight .....

```
mak@company $ ls proj
archive           footer.php      footer.php.latest
footer_final.php   header.php     header1.php
header2.php        header_current.php GodHelp
index.html         index.html.bkp  index.html.old
messed up         main_index.html main_header.php
never used        new_footer.php  new
old               old_data       todo
TODO.latest       toShowManager version1
version2          webHelp
:
:
```

[maktoons.blogspot.com/2009/06/if-dont-use-version-control-system.html](http://maktoons.blogspot.com/2009/06/if-dont-use-version-control-system.html)

## 2.3 git

We will use git, developed by Linus Torvalds, the “Linu” in Linux. In git each user stores a complete local copy of the project, including the history and all versions. So you do not rely as much on a centralized (remote) server. We will use bitbucket.org – it gives you unlimited free private repositories if you register with an academic email! First, install and configure git:

```
$ sudo apt-get install git
$ git config --global user.name "Your Name"
$ git config --global user.email "your.login@imperial.ac.uk"
$ git config --list
```

## 2.4 Your first repository

Time to bring your CMEECourseWork under version control:

```
$ cd CMEECourseWork
$ git init
$ echo "My CMEE 2014-14 Coursework Repository" > README.txt
$ git config --list
$ ls -all
$ git add README.txt #Staging
$ git status
$ git commit -m "Added README file." #you can use -am too
$ git status #what does it say now?
$ git add -A
$ git status
```

Nothing has been sent to a remote server yet (see section 2.5.1)! So let's go to bitbucket.org and setup:

- Login to your account
- Set up your ssh based access to bitbucket – <https://confluence.atlassian.com/bitbucket/set-up-ssh-for-git-728138079.html>
- Then create repository there with name CMEECourseWork
- Then grab the repository url and use `git remote add origin https... —check bitbucket help online.`

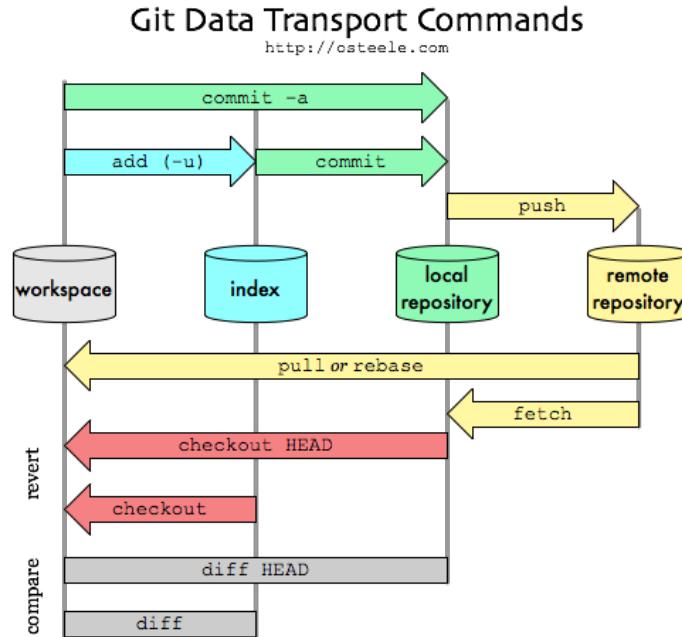
## 2.5 git commands

Here are some basic git commands:

<code>git init</code>	Initialize a new repository
<code>git clone</code>	Download a repository from a remote server
<code>git status</code>	Show the current status
<code>git diff</code>	Show differences between commits
<code>git blame</code>	Blame somebody for the changes!
<code>git log</code>	Show commit history
<code>git commit</code>	Commit changes to current branch
<code>git branch</code>	Show branches
<code>git branch name</code>	Create new branch
<code>git checkout name</code>	Switch to a different commit/branch
<code>git pull</code>	Upload from remote repository
<code>git push</code>	Send changes to remote repository

### 2.5.1 git command structure

Here is a graphical outline of the git command structure. Note that only when you `push` or `fetch` do you need an internet connection, as before that you are only archiving in a local (hidden) repository.



Keep in mind, the main mantra is, “commit often, comment always”!

	COMMENT	DATE
o	CREATED MAIN LOOP & TIMING CONTROL	14 HOURS AGO
o	ENABLED CONFIG FILE PARSING	9 HOURS AGO
o	MISC BUGFIXES	5 HOURS AGO
o	CODE ADDITIONS/EDITS	4 HOURS AGO
o	MORE CODE	4 HOURS AGO
o	HERE HAVE CODE	4 HOURS AGO
o	AAAAAAA	3 HOURS AGO
o	ADKFJSLKDFJSOKLFJ	3 HOURS AGO
o	MY HANDS ARE TYPING WORDS	2 HOURS AGO
o	HAAAAAAAAANDS	2 HOURS AGO

AS A PROJECT DRAGS ON, MY GIT COMMIT MESSAGES GET LESS AND LESS INFORMATIVE.

## 2.6 Ignoring Files

You will have some files you don't want to track (log files, temporary files, executables, etc). You can ignore entire classes of files with `.gitignore` (be in your CMEECourseWork!):

```
$ echo -e "*~\n*.tmp" > .gitignore
$ cat .gitignore
*~
*.tmp

$ git add .gitignore
$ touch temporary.tmp
$ git add *
The following paths are ignored by one of your .gitignore
files:
```

```
temporary.tmp
Use -f if you really want to add them.
fatal: no files added
```

## 2.7 Removing files

To remove a file use `git rm`:

```
$ echo "Text in a file to remove" > FileToRem.txt
$ git add FileToRem.txt
$ git commit -am "added a new file that we'll remove later"
master 5df9e96 added a new file that we'll remove later
 1 files changed, 1 insertions(+), 0 deletions(-)
 create mode 100644 FileToRem.txt

$ git rm FileToRem.txt
rm 'FileToRem.txt'

$ git commit -am "removed the file"
master b9f0bla removed the file
 1 files changed, 0 insertions(+), 1 deletions(-)
 delete mode 100644 FileToRem.txt
```

I typically just do all my stuff and then just use `git add -A`

## 2.8 Accessing history of the repository

To see particular changes introduced, read the repo's log :

```
$ git log
commit 08b5c1c78c8181d4606d37594681fdcfc3149ec
Author: Your Name <your.login@imperial.ac.uk>
Date:   Wed Oct 8 16:41:51 2014 -0500

    removed the file

commit 13f701775bce71998abe4dd1c48a4df8ed76c08b
Author: Your Name <your.login@imperial.ac.uk>
Date:   Wed Oct 5 16:41:16 2015 -0500

    added a new file that we'll remove later

commit a228dd3d5b1921ef18c5efd926ef11ca47306ed5
Author: Your Name <your.login@imperial.ac.uk>
Date:   Wed Oct 5 10:03:40 2015 -0500

    Added README file
```

For a more detailed version, add `-p` at the end.

## 2.9 Reverting to a previous version

If things horribly wrong with new changes, you can revert to the previous, “pristine” state:

```
$ git reset --hard
$ git commit -am "returned to previous state" #Note I used -am here
```

If instead you want to move back in time (temporarily), first find the “hash” for the commit you want to revert to, and then check-out:

```
$ git status
# On branch master
nothing to commit (working directory clean)

$ git log
commit c797824c9acbc59767a3931473aa3c53b6834aae
Author: Your Name <your.login@imperial.ac.uk>
Date:   Wed Aug 22 16:59:02 2014 -0500
.
.
.

$ git checkout c79782
```

Now you can play around. However, if you commit changes, you create a “branch” (git plays safe!). To go back to the future, type `git checkout master`

## 2.10 Branching

Imagine you want to try something out, but you’re not sure it will work well. For example, say you want to rewrite the Introduction of your paper, using a different angle, or you want to see whether switching to a library for a piece of code improves speed. What you then need is branching, which creates an “alternate reality” in which you can experiment:

```
$ git branch anexperiment

$ git branch
  anexperiment
* master

$ git checkout anexperiment
Switched to branch 'anexperiment'

$ git branch
* anexperiment
  master

$ echo "Do I like this better?" >> README.txt

$ git commit -am "Testing experimental branch"
[anexperiment 9f17dc1] Testing experimental branch
 1 files changed, 2 insertions(+), 0 deletions(-)
```

If you decide to merge the new branch after modifying it:

```
$ git checkout master

$ git merge anexperiment
Updating 08b5c1c..9f17dc1
Fast-forward
 README.txt |      2 ++
 1 files changed, 2 insertions(+), 0 deletions(-)

$ cat README.txt
My CMEE 2014-14 Coursework Repository
Do I like this better?
```

If there are no conflicts (i.e., some files that you changed also changed in the master in the meantime), you are done, and you can delete the branch:

```
$ git branch -d anexperiment
Deleted branch anexperiment (was 9f17dc1).
```

If instead you are not satisfied with the result, and you want to abandon branch:

```
$ git branch -D anexperiment
```

When you want to test something out, always branch! Reverting changes, especially in code, is typically painful.

## 2.11 Running git commands on a different directory

Since git version git 1.8.5, you can run git directly on a different directory than the current one using absolute or relative paths. For example, using a relative path, you can do:

```
git -C ../SomeDir/ status
```

## 2.12 Running git commands on multiple repositories at once

For git pulling in multiple subdirectories (each a separate repository):

```
$ find . -mindepth 1 -maxdepth 1 -type d -print -exec git -C {} pull \;
```

Breaking down these commands one by one,

find . searches the current directory  
 -type d to find directories, not files  
 -mindepth 1 for setting min search depth to one sub-directory  
 -maxdepth 1 for setting max search depth to one sub-directory

-exec git -C {} pull \ runs a custom git command for every git repo found

## 2.13 Wrapping up

- Invite me ([s.pawar@imperial.ac.uk](mailto:s.pawar@imperial.ac.uk)) to your CMEECourseWork repository
- CMEE2015MasteRepo will contain data and code files for upcoming practicals
- You will clone CMEE2015MasteRepo using `git clone git@bitbucket.org:mhasoba/cmee2015masterepo.git`
- You will thereafter `git pull CMEE2015MasteRepo`
- You will `git pull` inside CMEE2015MasteRepo thereafter (always use `git status` first)
- You will `cp` files from CMEE2015MasteRepo to your CMEECourseWork as and when needed

## 2.14 Practical

1. The only practical submission for git is the `.gitignore` and overall git repository `readme` file — make sure these in your coursework repository.

And of course, if you haven't gotten git with bitbucket going, you won't be able to submit any of your practicals anyway!

## 2.15 Readings & Resources

- Excellent book on Git: <http://git-scm.com/book>
- Also, <https://www.atlassian.com/git/>, including Bitbucket 101

# Chapter 3

## Advanced UNIX: Shell scripting

### 3.1 Shell scripting: What and Why

Instead of typing all the UNIX commands we need to perform one after the other, we can save them all in a file (a “script”) and execute them all at once.

The `bash` shell we are using provides a proper syntax that can be used to build complex command sequences and scripts.

Scripts can be used to automate repetitive tasks, to do simple data manipulation or to perform maintenance of your computer (e.g., backup). Indeed, most data manipulation can be handled by scripts without the need of writing a proper program.

### 3.2 Scripting: How

There are two ways of running a script, say `myscript.sh`:

1. The first is to call the interpreter `bash` to run the file

```
$ bash myscript.sh # OR sh myscript.sh
```

(A script that does something specific in a given project)

2. OR, make the script executable and execute it

```
$ chmod +x myscript.sh  
$ myscript.sh
```

(A script that does something generic, and is likely to be reused again and again – can you think of examples?)

The generic scripts of type (2) can be saved in `username/bin/` and made executable (the `.sh` extension not needed)

```
$ mkdir ~/bin
$ PATH=$PATH:$HOME/bin #Tell UNIX to look in /home/bin for commands
```

### 3.3 Your first shell script

Let's write our first shell script! For starters,

- ★ Write and save `boilerplate.sh` in `CMEECourseWork/Week1/Code`, and add the following script to it:

```
#!/bin/bash
# Author: Your Name your.login@imperial.ac.uk
# Script: boilerplate.sh
# Desc: simple boilerplate for shell scripts
# Arguments: none
# Date: Oct 2015

echo -e "\nThis is a shell script! \n" #what does -e do?

#exit
```

The first line is a “shebang” (or sha-bang or hashbang or pound-bang or hash-exclam or hashpling! – Wikipedia). It can also can be written as `#!/bin/sh`. It tells the bash interpreter that this is a bash script and that it should be interpreted and run as such. The hash marks in the following lines tell the interpreter that it should ignore the lines following them (that's how you put in script documentation (who wrote the script and when, what the script does, etc.) and comments on particular line of script.

Now run the script by typing in the terminal:

```
$ bash boilerplate.sh
```

### 3.4 A useful shell-scripting example

Let's write a shell script to transform comma-separated files (csv) to tab-separated files and vice-versa. This can be handy — for example, in certain computer languages, it is much easier to read tab or space separated files than csv (e.g., C)

To do this, in the bash we can use `tr`, which deletes or substitute characters. Here are some examples.

```
$ echo "Remove      excess      spaces." | tr -s "\b" " "
Remove excess spaces.
$ echo "remove all the as" | tr -d "a"
remove ll the s
$ echo "set to uppercase" | tr [:lower:] [:upper:]
SET TO UPPERCASE
$ echo "10.00 only numbers 1.33" | tr -d [:alpha:] |
  tr -s "\b" ", "
```

```
10.00,1.33
```

Now write a shell script to substitute all tabs with commas called `tabtocsv.sh` in Week1/Code:

```
#!/bin/bash
# Author: Your name you.login@imperial.ac.uk
# Script: tabtocsv.sh
# Desc: substitute the tabs in the files with commas
#       saves the output into a .csv file
# Arguments: 1-> tab delimited file
# Date: Oct 2015

echo "Creating a comma delimited version of $1 ..."

cat $1 | tr -s "\t" "," >> $1.csv

echo "Done!"

exit
```

Now test it (note where the output file gets saved)

```
echo -e "test \t\t test" >> ../SandBox/test.txt
bash tabtocsv.sh ../SandBox/test.txt
```

## 3.5 Variables in shell scripting

There are three ways to assign values to variables (note lack of spaces!):

1. Explicit declaration: `MYVAR=myvalue`
2. Reading from the user: `read MYVAR`
3. Command substitution: `MYVAR=$( (ls | wc -l) )`

Here are some examples of assignments (try it out save as Week1/Code/variables.sh):

```
#!/bin/bash
# Shows the use of variables
MyVar='some string'
echo 'the current value of the variable is' $MyVar
echo 'Please enter a new string'
read MyVar
echo 'the current value of the variable is' $MyVar
## Reading multiple values
echo 'Enter two numbers separated by space(s)'
read a b
echo 'you entered' $a 'and' $b '. Their sum is:'
mysum=`expr $a + $b`
echo $mysum
```

And also (save as Week1/Code/MyExampleScript.sh):

```
#!/bin/bash

msg1="Hello"
msg2=$USER
echo "$msg1 $msg2"

echo "Hello $USER"
echo
```

### 3.6 Some more Examples

Here are a few more illustrative examples (test each one out, save in Week1/Code/ with the given name):

CountLines.sh:

```
#!/bin/bash
NumLines=`wc -l < $1`
echo "The file $1 has $NumLines lines"
echo
```

ConcatenateTwoFiles.sh:

```
#!/bin/bash
cat $1 > $3
cat $2 >> $3
echo "Merged File is"
cat $3
```

### 3.7 Practical

#### 1. Some instructions:

Along with the completeness of the practicals/exercises themselves, you will be marked on the basis of how complete and well-organized your directory structure and content is.

Review (especially if you got lost along the way) and make sure all your shell scripts are functional: boilerplate.sh, ConcatenateTwoFiles.sh, CountLines.sh, MyExampleScript.sh, tabtocsv.sh, variables.sh

Don't worry about how some of these scripts will run on my computer without explicit inputs (e.g., ConcatenateTwoFiles.sh needs two input files) — I will run them with my own test files.

Make sure you have your weekly directory organized with Data, Sandbox, Code with the necessary files, under CMEECourseWork/Week1. *All scripts should run on any other Unix/Linux machine* — for example, always call data from the Data directory using relative paths.

**Make sure there is a `readme` file in every week's directory. This file should give an overview of the weekly directory contents, listing all the scripts and what they do. This is different from the `readme` for your overall git repository, of which Week 1 is a part. You will write a similar `readme` for each subsequent weekly submission.**

**Don't put any scripts that are part of the submission in your `home/bin` directory! You can put a copy there, but a working version should be in your repository.**

2. Finally, a small exercise: write a `csvtospace.sh` shell script that takes a comma separated values and converts it to a space separated values file. However, it must not change the input file — it should save it as a differently named file.

Save the script in `CMECourseWork/Week1/Code`, and run it on the `csv` data files that are in `Temperatures` in the master repository's `Data` directory.

*Don't modify anything (or refer to anything) in your local copy of the master repository. All changes you make in the master repository will be lost. Copy whatever you need from the master repository to your own repository.*

*Commit and push everything by next Wednesday 5 PM.*

This includes `UnixPrac1.txt`! Check the updated instructions from Chapter 1 on this practical.

## 3.8 Readings & Resources

- Plenty of shell scripting resources and tutorials out there; in particular, look up <http://www.tutorialspoint.com/unix/unix-using-variables.htm>



## Chapter 4

# Using L<sup>A</sup>T<sub>E</sub>X for scientific documents

### 4.1 What's L<sup>A</sup>T<sub>E</sub>X?

In your research, you will produce papers, reports and – very importantly – your thesis. These documents can be written using a WYSIWYG (What You See Is What You Get) editor (e.g., Word). However, an alternative especially suited for scientific publications is L<sup>A</sup>T<sub>E</sub>X. In L<sup>A</sup>T<sub>E</sub>X, the document is simply a text file (`.tex`). Text formatting is using markups (like HTML). The file is then “compiled” (like source code of a programming language) into a file – typically `.pdf`.

### 4.2 Why L<sup>A</sup>T<sub>E</sub>X?

- The input is a small, portable text file
- L<sup>A</sup>T<sub>E</sub>X compilers are freely available for all OS'
- Exactly the same result on any computer (not true for Word)
- L<sup>A</sup>T<sub>E</sub>X produces beautiful, professional looking docs (e.g., like this one!)
- Mathematical formulas (esp complex ones) are easy to write
- L<sup>A</sup>T<sub>E</sub>X is very stable – current version basically same since 1994! (9 major versions of MS Word since 1994 – with compatibility issues)
- L<sup>A</sup>T<sub>E</sub>X is free!
- Many journals provide L<sup>A</sup>T<sub>E</sub>X templates, making formatting quicker
- Bibliographies are a breeze and work with Mendeley and Zotero
- Plenty of online support available – your question has probably already been answered
- You can integrate L<sup>A</sup>T<sub>E</sub>X into a workflow to auto-generate lengthy and complex documents (like your thesis).



### 4.2.1 Limitations of LATEX

- It has a steeper learning curve.
- Can be difficult to manage revisions with multiple authors – especially if they don't use LATEX! (I have a dark secret)
- Typesetting tables can be a bit complex.
- Images and floats don't jump like Word, but if you don't use the right package, they can be difficult to place where you want!

## 4.3 Installing LATEX

```
sudo apt-get install texlive-full texlive-fonts-recommended
latex-beamer texpower texlive-pictures texlive-latex-extra
texpower-examples imagemagick
```

We will use a text editor in this lecture, but you can use one of a number of WYSIWYG frontends (e.g., Lyx, TeXmacs), as well as GUI's (texmaker, Gummi, TeXShop, etc)

## 4.4 A first LATEX example

- ★ Open geany and type the following in a file Week1/Code/FirstExample.tex:

```
\documentclass[12pt]{article}
\title{A Simple Document}
\author{Your Name}
\date{}
\begin{document}
\maketitle

\begin{abstract}
This paper must be cool!
\end{abstract}

\section{Introduction}
But I didn't understand anything.

\section{Materials \& Methods}
One of the most famous equations is:
\begin{equation}
E = mc^2
\end{equation}
This equation was first proposed by Einstein in 1905
\cite{einstein1905does}.

\bibliographystyle{plain}
\bibliography{FirstBiblio}
\end{document}
```

Now, let's get a citation for Einstein's paper:

- \* In Google Scholar, go to “settings” (upper right corner) and choose BibTeX as bibliography manager.
- \* Now type “does the energy of a body einstein 1905”
- \* The paper should be the one on the top.
- \* Click “Import into BibTeX” should show the following text, that you will save in the file FirstBiblio.bib (in the same directory as FirstExample.tex)

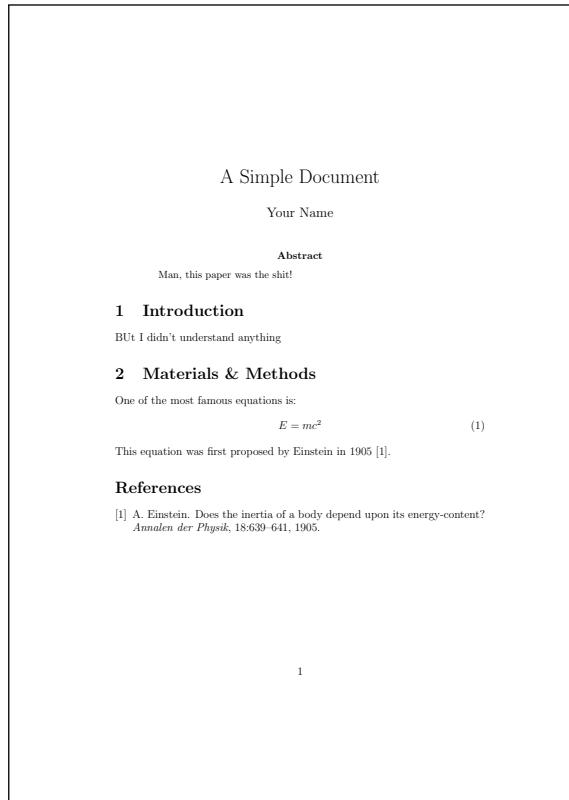
```
@article{einstein1905does,
    title={Does the inertia of a body depend upon its energy-content?},
    author={Einstein, A.},
    journal={Annalen der Physik},
    volume={18},
    pages={639--641},
    year={1905}
}
```

Now we can create a .pdf of the article.

- \* In the terminal type (are you the right directory?!):

```
$ pdflatex FirstExample.tex
$ pdflatex FirstExample.tex
$ bibtex FirstExample
$ pdflatex FirstExample.tex
$ pdflatex FirstExample.tex
```

This should produce the file FirstExample.pdf:



#### 4.4.1 A bash script to compile LATEX

You can of course write a useful little bash script to compile latex with bibtex!

Type the following script and call it `CompileLaTeX.sh` (you know where to put it!):

```
#!/bin/bash
pdflatex $1.tex
pdflatex $1.tex
bibtex $1
pdflatex $1.tex
pdflatex $1.tex
evince $1.pdf &

## Cleanup
rm *~
rm *.aux
rm *.dvi
rm *.log
rm *.nav
rm *.out
rm *.snm
rm *.toc
```

How do you run this script? The same as your previous bash scripts, so

```
$ bash CompileLaTeX.sh FirstExample
```

Why have I not written the `*.tex` extension of `FirstExample` in the command above?

## 4.5 A brief LATEX tour

- Spaces, new lines and special characters:
  - Several spaces in your text editor are treated as one space in the typeset document
  - Several empty lines are treated as one empty line
  - One empty line defines a new paragraph
  - Some characters are “special”: # \$ % ^ & \_ { } ~ \.
  - To type these special characters, you have to add a “backslash” in front, e.g., \\$ produces \$.
- Document structure:
  - Each LATEX command starts with \ (e.g., to get LATEX, you need \LaTeX)
  - The first command is always \documentclass defining the type of document (e.g., article, book, report, letter).
  - You can set several options. For example, to set size of text to 10 points and the letter paper size: \documentclass[10pt,letterpaper]{article}.
- After having declared the type of document, you can specify packages you want to use. The most useful are:
  - \usepackage{color}: use colors for text in your document.

- `\usepackage{amsmath, amssymb}`: American Mathematical Society formats and commands for typesetting mathematics.
  - `\usepackage{fancyhdr}`: fancy headers and footers.
  - `\usepackage{graphicx}`: include figures in pdf, ps, eps, gif and jpeg.
  - `\usepackage{listings}`: typeset source code for various programming languages.
  - `\usepackage{rotating}`: rotate tables and figures.
  - `\usepackage{lineno}`: line numbers.
- Once you select the packages, you can start your document with `\begin{document}`, and end it with `\end{document}`.

## 4.6 LATEX templates

There are lots of useful LATEX templates out there. As an example of structure of a document, take the article template provided by the journal PNAS:

```
\documentclass{pnastwo}
\usepackage{amssymb, amsfonts, amsmath}
%% For PNAS Only:
\contributor{Submitted to Proceedings
of the National Academy of Sciences of the United States of America}
\url{www.pnas.org/cgi/doi/10.1073/pnas.0709640104}
\copyrightyear{2014}
\issuedate{Issue Date}
\volume{Volume}
\issuenumber{Issue Number}

\begin{document}
\title{My Title}
\author{Some Name \affil{1}{Imperial College London, UK} \and
Some O. Name\affil{2}{University of Exeter, Penryn, Cornwall, UK}}
\maketitle
\begin{article}
\begin{abstract}
Mind blowing abstract.
\end{abstract}
\keywords{term1 | term2 | term3}

%% Main text of the paper
\dropcap{I}n this work, we show how \LaTeX can be used to typeset a PNAS ↔
paper. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Phasellus↔
sodales consectetur lobortis. Proin tincidunt eros dapibus ipsum ↔
faucibus sed rhoncus augue mollis. In lectus velit, interdum at ↔
adipiscing quis, imperdiet sed justo. Praesent commodo, mi iaculis ↔
tincidunt mollis, sapien lectus aliquam neque, ac faucibus arcu est eu ↔
sem. Ut non lacus lacus, eu suscipit odio. Aliquam erat volutpat. Vivamus↔
dapibus pretium nunc, et placerat turpis bibendum mollis. Fusce eu mi ut↔
nulla accumsan viverra. In nulla tellus, ultrices ut venenatis nec, ↔
laoreet eget diam. Pellentesque aliquam facilisis ultricies. Vestibulum ↔
sollicitudin leo non neque vehicula a volutpat eros faucibus. Vestibulum ↔
nec lorem dui.

\begin{materials}
These are the materials and methods.
\end{materials}
```

```
\begin{acknowledgments}
-- text of acknowledgments here, including grant info --
\end{acknowledgments}

\end{article}
\end{document}
```

I have added some templates in the CMEEMasteRepo that you should have a look and play around with

## 4.7 Typesetting math

There are two ways to display math

1. First, one can produce inline mathematics (i.e., within the text).
2. Second, one can produce stand-alone, numbered equations and formulae.

For inline math, the “dollar” sign flanks the math to be typeset. For example, the code:

```
$\int_0^1 p^x (1-p)^y dp$
```

Becomes  $\int_0^1 p^x (1-p)^y dp$

For numbered equations (almost always a great idea), LATEX provides the `equation` environment:

```
\begin{equation}
\int_0^1 \left( \ln \left( \frac{1}{x} \right) \right)^y dx = y!
\end{equation}
```

Becomes:

$$\int_0^1 \left( \ln \left( \frac{1}{x} \right) \right)^y dx = y! \quad (4.1)$$

## 4.8 A few more tips

The following tips might prove handy:

- LATEX has a full set of symbols and operators (plenty of lists online)
- Long documents can be split into separate `.tex` documents and combined using `input`
- Long documents can be split into separate `.tex` documents and
- Figures can be included using the `graphicx` package
- You can use Mendeley to export and maintain `.bib` files
- You can redefine environments and commands in the preamble

## 4.9 Readings & Resources

- The Visual  $\text{\LaTeX}$  FAQ: sometimes it is difficult to describe what you want to do!  
<http://get-software.net/info/visualFAQ/visualFAQ.pdf>
- Myriad online resources for  $\text{\LaTeX}$ , including:  
[www.http://en.wikibooks.org/wiki/LaTeX/Introduction](http://en.wikibooks.org/wiki/LaTeX/Introduction),  
[www.ctan.org/tex-archive/info/lshort/english/](http://www.ctan.org/tex-archive/info/lshort/english/)  
<http://ftp.uni-erlangen.de/mirrors/CTAN/info/lshort/english/lshort.pdf>
- Beautiful presentations in  $\text{\LaTeX}$ : <http://tug.org/pracjourn/2005-2/miller/miller.pdf>
- Bibliographies in  $\text{\LaTeX}$ : <http://schneider.ncifcrf.gov/latex.html>

## 4.10 Practical: $\text{\LaTeX}$ , and wrap up the week's work

Test `CompileLaTeX.sh` with `FirstExample.tex` and bring it under version control under `CMEECourseWork/Week1` in your repository. Make sure that `CompileLaTeX.sh` will work if I ran it from my computer using `FirstExample.tex` as an input.

Again, make sure you have your Week 1 directory organized with Data, Sandbox and Code with the necessary files and this week's (functional!) scripts in there. Every script should run without errors on my computer. This includes the five solutions (single-line commands you came up with) in `UnixPrac1.txt`!

*Commit and push everything by next Wednesday 5 PM.*



# Chapter 5

## Biological Computing in Python I

Science is what we understand well enough to explain to a computer. Art is everything else we do

---

—Donald Knuth

### 5.1 Chapter 1's UNIX question?

```
find . -type f -exec ls -s {} \; | sort -n | head -10
```

### 5.2 Outline of the python module

The content and structure of the Python weeks are geared towards teaching you:

- Scientific programming in biology using `python`
- Basics of `python`
- To write and run `python` code, understand and implement control flows
- Learn to use `ipython`
- Writing, debugging, using, and testing `python` functions
- Efficient numerical programming in `python`
- Regular expressions, certain `python` packages, `python` with databases
- Using `python` to run other “stuff”, and to patch together data analyses or numerical simulation workflows

### 5.3 Why `python`?

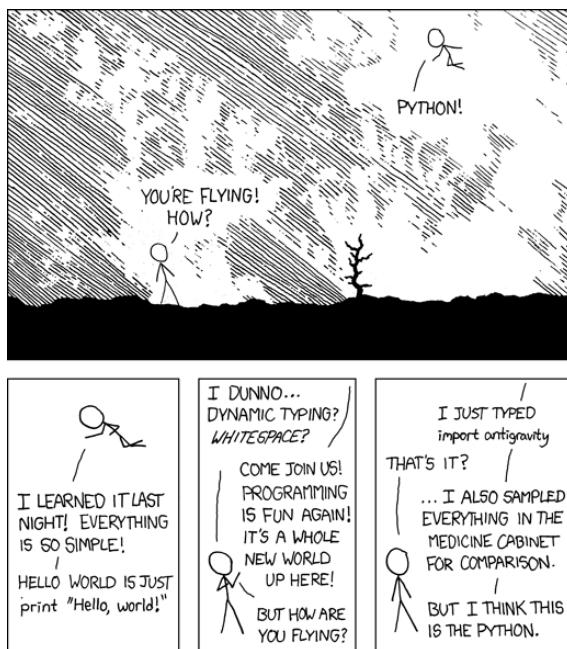
`python` was designed with readability and re-usability in mind. Time taken by programming + debugging + running is likely to be relatively lower in `python` than less intuitive or cluttered languages (e.g., FORTRAN, `perl`). It is a pretty good solution if you want to easily write readable code that is also reasonably efficient (computationally speaking).

	<b>Fortran</b>	<b>Julia</b>	<b>Python</b>	<b>R</b>	<b>Matlab</b>	<b>Octave</b>	<b>Mathematica</b>	<b>JavaScript</b>	<b>Go</b>	<b>LuaJIT</b>	<b>Java</b>
	gcc 5.1.1	0.4.0	3.4.3	3.2.2	R2015b	4.0.0	10.2.0	V8 3.28.71.19	go1.5	gsl-shell 2.3.1	1.8.0_45
fib	0.70	2.11	77.76	533.52	26.89	9324.35	118.53	3.36	1.86	1.71	1.21
parse_int	5.05	1.45	17.02	45.73	802.52	9581.44	15.02	6.06	1.20	5.77	3.35
quicksort	1.31	1.15	32.89	264.54	4.92	1866.01	43.23	2.70	1.29	2.03	2.60
mandel	0.81	0.79	15.32	53.16	7.58	451.81	5.13	0.66	1.11	0.67	1.35
pi_sum	1.00	1.00	21.99	9.56	1.00	299.31	1.69	1.01	1.00	1.00	1.00
rand_mat_stat	1.45	1.66	17.93	14.56	14.52	30.93	5.95	2.30	2.96	3.27	3.92
rand_mat_mul	3.48	1.02	1.14	1.57	1.12	1.12	1.30	15.07	1.42	1.16	2.36

**Figure:** benchmark times relative to C (smaller is better, C performance = 1.0).

<http://julialang.org/>

Figure 5.1: python is pretty fast!



[www.xkcd.com](http://www.xkcd.com)

Figure 5.2: Is python the most common answer to your daily programming needs? Possibly!

### 5.3.1 The Zen of python

Open a terminal and type

```
$ python -c "import this"
```

## 5.4 Installing python

We will use 2.7.x, not 3.x (you can use 3.x later, if you want)

On Ubuntu/Linux, open a terminal (ctrl+alt+t) and type:

```
sudo apt-get update
sudo apt-get install python
```

```
sudo apt-get install ipython python-scipy python-matplotlib
```

**Note:** Geany users can enable “Send Selection to Terminal” using <Primary>Return by going to Edit > Preferences > Keybindings.

## 5.5 python warmup

Open a terminal (ctrl+alt+t) and type python (or use the terminal that you just used to install python). Then, try the following:

```
>>> 2 + 2 # Summation; note that comments start with #
4

>>> 2 * 2 # Multiplication
4

>>> 2 / 2 # Integer division
1

>>> 2 / 2.0 # "Float" division, note the output is float
1.0

>>> 2 / 2.
1.0

>>> 2 > 3
False

>>> 2 >= 2
True
```

What does “float” mean in the above comment? Why is it necessary to specify this in Python (not necessary in Python 3.x)?

Now let’s switch to a python interface called ipython. We will see later why. Type ctrl+D: this will exit you from the python shell and you will see the bash prompt again. Then type ipython. You should now see (after some text):

```
In [1]:
```

This is the interactive python shell (or, “ipython”). If you don’t like the blue prompt, type %colors linux

### 5.5.1 Variable types

Now, let’s continue our warmup (ignore the prompt numbering [1], [2], etc)

```
In []: 2 == 2
Out []: True
```

```
In []: 2 != 2
Out []: False

In []: 3 / 2
Out []: 1

In []: 3 / 2.
Out []: 1.5

In []: 'hola, ' + 'mi llamo Samraat' #why not two languages at the same
time?!
Out []: 'hola, mi llamo Samraat'
```

Thus, python has integer, float (real numbers, with different precision levels) and string variables.

### 5.5.2 python operators

Here are are the operators in python:

+	Addition
-	Subtraction
*	Multiplication
/	Division
**	Power
%	Modulo
//	Integer division
==	Equals
!=	Differs
>	Greater
>=	Greater or equal
&, and	Logical and
, or	Logical or
!, not	Logical not

### 5.5.3 Assigning and manipulating variables

```
In []: x = 5

In []: x + 3
Out []: 8

In []: y = 8

In []: x + y
Out []: 13

In []: x = 'My string'

In []: x + ' now has more stuff'
Out []: 'My string now has more stuff'
```

```
In []: x + y
Out []: TypeError: cannot concatenate 'str' and 'int' objects
```

OK, so concatenating string and numeric (integer in this case) variables doesn't work. No problem, we can convert from one type to another:

```
In []: x + str(y)
Out []: 'My string8'

In []: z = '88'

In []: x + z
Out []: 'My string88'

In []: y + int(z)
Out []: 96
```

Note that in python, the type of a variable is determined when the program or command is running (dynamic typing) (like R, unlike C or FORTRAN).

## 5.6 python data types and data structures

python number or string variables (or both) can be stored and manipulated in:

- **List**: most versatile, can contain compound data, “mutable”, enclosed in brackets, [ ]
- **Tuple**: like a list, but “immutable” — like a read only list, enclosed in parentheses, ()
- **Dictionary**: a kind of “hash table” of key-value pairs enclosed by curly braces, { } — key can be number or string, values can be any object! (well OK, a python object)
- **numpy arrays**: Fast, compact, convenient for numerical computing — more on this later!

### 5.6.1 Lists

```
In []: myList = [3, 2.44, 'green', True]

In []: myList[1]
Out []: 2.44

In []: myList[0] # NOTE: FIRST ELEMENT -> 0
Out []: 3

In []: myList[4]
Out []: IndexError: list index out of range

In []: myList[2] = 'blue'

In []: myList
Out []: [3, 2.44, 'blue', True]

In []: myList[0] = 'blue'

In []: myList
```

```

Out []: ['blue', 2.44, 'blue', True]

In []: myList.append('a new item') # NOTE: ".append"!

In []:
Out []: ['blue', 2.44, 'blue', True, 'a new item']

In []: myList.sort() # NOTE: suffix a ".", hit tab, and wonder!

In []:
Out []: [True, 2.44, 'a new item', 'blue', 'blue']

```

In the above commands, notice that python “indexing” starts at 0, not 1!

### 5.6.2 Tuples

```

In []: FoodWeb=[('a', 'b'), ('a', 'c'), ('b', 'c'), ('c', 'c')]

In []:
Out []: ('a', 'b')

In []:
Out []: 'a'

In []:
FoodWeb[0][0] = "bbb" # NOTE: tuples are "immutable"
TypeError: 'tuple' object does not support item assignment

In []:
FoodWeb[0] = ("bbb", "ccc")

In []:
Out []: ('bbb', 'ccc')

```

Note that tuples are “immutable”; that is, a particular pair or sequence of strings or numbers cannot be modified after it is created.

In the above example, why assign these food web data to a list of tuples and not a list of lists? — because we want to maintain the species associations, no matter what — they are sacrosanct!

Tuples contain immutable sequences, but you can append to them:

```

In []: a = (1, 2, [])
In []:
a[2].append(1000)

In []:
Out []: (1, 2, [1000])

```

### 5.6.3 Sets

You can convert a list to an immutable “set”— an unordered collection with no duplicate elements. Once you create a set you can perform set operations on it:

```
In []: a = [5,6,7,7,7,8,9,9]
In []: b = set(a)
In []: b
Out []: set([8, 9, 5, 6, 7])
In []: c = set([3,4,5,6])
In []: b & c
Out []: set([5, 6])
In []: b | c
Out []: set([3, 4, 5, 6, 7, 8, 9])
In []: list(b | c) # set to list
Out []: [3, 4, 5, 6, 7, 8, 9]
```

The key set operations in python are:

a - b	a.difference(b)
a <= b	a.issubset(b)
a >= b	b.issubset(a)
a & b	a.intersection(b)
a   b	a.union(b)

#### 5.6.4 Dictionaries

A set of values (any python object) indexed by keys (string or number), a bit like R lists.

```
In []: GenomeSize = {'Homo sapiens': 3200.0, 'Escherichia coli': 4.6,
'Arabidopsis thaliana': 157.0}

In []: GenomeSize
Out []:
{'Arabidopsis thaliana': 157.0,
 'Escherichia coli': 4.6,
 'Homo sapiens': 3200.0}

In []: GenomeSize['Arabidopsis thaliana']
Out []: 157.0

In []: GenomeSize['Saccharomyces cerevisiae'] = 12.1

In []: GenomeSize
Out []:
{'Arabidopsis thaliana': 157.0,
 'Escherichia coli': 4.6,
 'Homo sapiens': 3200.0,
 'Saccharomyces cerevisiae': 12.1}

In []: GenomeSize['Escherichia coli'] = 4.6 # ALREADY IN DICTIONARY!

In []: GenomeSize
Out []:
{'Arabidopsis thaliana': 157.0,
```

```
'Escherichia coli': 4.6,
'Homo sapiens': 3200.0,
'Saccharomyces cerevisiae': 12.1}

In []: GenomeSize['Homo sapiens'] = 3201.1

In []: GenomeSize
Out []:
{'Arabidopsis thaliana': 157.0,
 'Escherichia coli': 4.6,
 'Homo sapiens': 3201.1,
 'Saccharomyces cerevisiae': 12.1}
```

So, in summary,

- If your elements/data are unordered and indexed by numbers use **lists**
- If they are ordered sequences use a **tuple**
- If you want to perform set operations on them, use a **set**
- If they are unordered and indexed by keys (e.g., names), use a **dictionary**

*But why not use dictionaries for everything?* – because it can slow down your code!

### 5.6.5 Copying mutable objects

Copying mutable objects can be tricky. Try this:

```
# First, try this:
a = [1, 2, 3]
b = a # you are merely creating a new "tag" (b)
a.append(4)
print b
# this will print [1, 2, 3, 4]!!

# Now, try:
a = [1, 2, 3]
b = a[:] # This is a "shallow" copy
a.append(4)
print b
# this will print [1, 2, 3].

# What about more complex lists?
a = [[1, 2], [3, 4]]
b = a[:]
a[0][1] = 22 # Note how I accessed this 2D list
print b
# this will print [[1, 22], [3, 4]]

# the solution is to do a "deep" copy:
import copy

a = [[1, 2], [3, 4]]
b = copy.deepcopy(a)
a[0][1] = 22
print b
```

```
# this will print [[1, 2], [3, 4]]
```

So, you need to employ `deepcopy` to really copy an existing object or variable and assign a new name to the copy.

### 5.6.6 python with strings

One of the things that makes python so useful and versatile, is that it has a powerful set of inbuilt commands to perform string manipulations. For example, try these:

```
s = " this is a string "
len(s)
# length of s -> 18

print s.replace(" ", "-")
# Substitute spaces " " with dashes -> -this-is-a-string-

print s.find("s")
# First occurrence of s -> 4 (start at 0)

print s.count("s")
# Count the number of "s" -> 3

t = s.split()
print t
# Split the string using spaces and make
# a list -> ['this', 'is', 'a', 'string']

t = s.split(" is ")
print t
# Split the string using " is " and make
# a list -> [' this', 'a string ']

t = s.strip()
print t
# remove trailing spaces

print s.upper()
# ' THIS IS A STRING '

'WORD'.lower()
# 'word'
```

## 5.7 python Input/Output

Let's look at importing and exporting data. Make a textfile called `test.txt` in `Week2/Sandbox/` with the following content (including the empty lines):

```
First Line
Second Line

Third Line
```

Fourth Line
-------------

Then, type the following in Week2/Code/basic\_io.py:

```
#####
# FILE INPUT
#####
# Open a file for reading
f = open('../Sandbox/test.txt', 'r')
# use "implicit" for loop:
# if the object is a file, python will cycle over lines
for line in f:
    print line, # the "," prevents adding a new line

# close the file
f.close()

# Same example, skip blank lines
f = open('../Sandbox/test.txt', 'r')
for line in f:
    if len(line.strip()) > 0:
        print line,
f.close()

#####
# FILE OUTPUT
#####
# Save the elements of a list to a file
list_to_save = range(100)

f = open('../Sandbox/testout.txt', 'w')
for i in list_to_save:
    f.write(str(i) + '\n') ## Add a new line at the end

f.close()

#####
# STORING OBJECTS
#####
# To save an object (even complex) for later use
my_dictionary = {"a key": 10, "another key": 11}

import pickle

f = open('../Sandbox/testp.p', 'wb') ## note the b: accept binary files
pickle.dump(my_dictionary, f)
f.close()

## Load the data again
f = open('../Sandbox/testp.p', 'rb')
another_dictionary = pickle.load(f)
f.close()

print another_dictionary
```

Note the following:

The `csv` package makes it easy to manipulate CSV files (get `testcsv.csv` from CMEMasterRepo). Type the following script in Week2/Code/basic\_csv.py

```
import csv

# Read a file containing:
# 'Species','Infraorder','Family','Distribution','Body mass male (Kg)'
f = open('../Sandbox/testcsv.csv','rb')

csvread = csv.reader(f)
temp = []
for row in csvread:
    temp.append(tuple(row))
    print row
    print "The species is", row[0]

f.close()

# write a file containing only species name and Body mass
f = open('../Sandbox/testcsv.csv','rb')
g = open('../Sandbox/bodymass.csv','wb')

csvread = csv.reader(f)
csvwrite = csv.writer(g)
for row in csvread:
    print row
    csvwrite.writerow([row[0], row[4]])

f.close()
g.close()
```

## 5.8 Writing python code

Now let's learn to write and run python code from a `*.py` file. But first, some some guidelines for good code-writing practices (see [python.org/dev/peps/pep-0008/](http://python.org/dev/peps/pep-0008/)):

- Wrap lines to be <80 characters long. You can use parentheses () or signal that the line continues using a “backslash” \
- Use either 4 spaces for indentation or tabs, but not both! (I use tabs!)
- Separate functions using a blank line
- When possible, write comments on separate lines

Make sure you have chosen a particular indent type (space or tab) in geany (or whatever IDE you are using) — indentation is all-important in python. Furthermore,

- Use “docstrings” to **document how to use the code**, and **comments to explain why and how the code works**
- Naming conventions (bit of a mess, you'll learn as you go!):
  - `_internal_global_variable` (for use inside module only)
  - `a_variable`
  - `SOME_CONSTANT`

- a\_function
- Never call a variable l or O or o  
*why not?* – you are likely to confuse it with 1 or 0!
- Use spaces around operators and after commas:  
a = func(x, y) + other(3, 4)

### 5.8.1 Writing python functions (or modules)

Let's start with a “boilerplate” code. Type the code below and save as boilerplate.py in CMEECourseWork/Week2/Code:

```
#!/usr/bin/python

"""Description of this program
you can use several lines"""

__author__ = 'Samraat Pawar (s.pawar@imperial.ac.uk)'
__version__ = '0.0.1'

# imports
import sys # module to interface our program with the operating system

# constants can go here

# functions can go here
def main(argv):
    print 'This is a boilerplate' # NOTE: indented using two tabs or 4 ←
        spaces
    return 0

if (__name__ == "__main__"): #makes sure the "main" function is called from ←
    commandline
    status = main(sys.argv)
    sys.exit(status)
```

No cd to the directory and run the code:

```
$ cd ~/Documents/.../CMEECourseWork/Week2/Code
$ python boilerplate.py
```

You should see “This is a boilerplate” in your terminal window.

Alternatively, you can use ipython:

```
$ ipython boilerplate.py
```

### 5.8.2 Components of the python function

Now let's look at the elements of your first, boilerplate code:

- First “shebang” line tells computer where to look for python (similar to unix shell scripts)
- Triple quotes start a “docstring” comment (more on this later)
- “`__`” signal “internal” variables (never name your variables so!)
- “`def`” indicates the start of a python function; all subsequent lines will be indented
- Last few lines, including the `main` are esoteric but important; more on this below

### Python has “modules”

Firstly, it’s important to know that somewhat confusingly, Pythonistas call a file containing function definitions and statements (e.g., assignments of constant variables) a “module”. There is a practical reason (there’s always one!) for this. You might want to use a particular set of python def’s (functions) and statements either as a standalone function, or use it or subsets of it from other scripts. So in theory, every function you define can be a sub-module usable by other scripts.

*In other words, definitions from a module can be imported into other modules and scripts, or into the main module itself.*

At this juncture, you might also want to know that Python also has a “class” construct. Have a look at <http://learnpythonthehardway.org/book/ex40.html> — a nice, intuitive tutorial that should help you understand functions vs. modules vs. classes in Python.

### Why include `__name__ == "__main__"` and all that jazz

When you run a Python module with or without arguments, the code in the called module will be executed just as if you imported it, but with the `__name__` set to `"__main__"`. So adding this code at the end of your module,

```
if __name__ == "__main__":
```

directs the python interpreter to set the special `__name__` variable to have a value `"__main__"`, so that the file is usable as a script as well as an importable module. How do you import? Simply as (in python or ipython shell):

```
In []: import boilerplate
```

Then type

```
In []: boilerplate
Out[]: <module 'boilerplate' from 'boilerplate.py'>
```

One more script to hopefully clarify this further. Type and save the following in a script file called `using_name.py`:

```
#!/usr/bin/python
# Filename: using_name.py

if __name__ == '__main__':
```

```

    print 'This program is being run by itself'
else:
    print 'I am being imported from another module'

```

Now run it:

```

In []: run using_name.py
This program is being run by itself

```

Now, try:

```

In []: import using_name
I am being imported from another module

```

The output I am being imported from another module will only show up once.

Also please look up <https://docs.python.org/2/tutorial/modules.html>

### What on earth is sys.argv?

In your boilerplate code, as any other Python code, argv is the “argument variable.” Such variables are necessarily very common across programming languages, and play an important role — argv is a variable that holds the arguments you pass to your Python script when you run it. sys.argv is simply the an object created by python using the sys module (which you imported at the beginning of the script) that contains the names of the argument variables in the current script.

To understand this in a practical way, let’s write and save a script called sysargv.py:

```

import sys
print "This is the name of the script: ", sys.argv[0]
print "Number of arguments: ", len(sys.argv)
print "The arguments are: " , str(sys.argv)

```

Now run sysargv.py with different numbers of arguments:

```

run sysargv.py
run sysargv.py var1 var2
run sysargv.py 1 2 var3

```

As you can see the first variable is always the file name, and is always available as to the Python interpreter.

Then, the command main(argv=sys.argv) directs the interpreter to pass the argument variables to the main function. Which brings us to,

```

def main(argv):
    print 'This is a boilerplate' # NOTE: indented using two tabs or
        four spaces

```

This is the main function. Arguments obtained in the `if __name__ == "__main__":` part of the script are “fed” to this main function where the printing of the line “This is a boilerplate” happens.

OK, finally, what about this bit:

```
sys.exit(status)
```

It's just a way to exit in a dignified (or abrupt) manner (from anywhere in the program)! Try putting it elsewhere in a function/module and see what happens.

### 5.8.3 Variable scope

One important thing to note about functions, in any language, is that variables inside functions are invisible outside of it, nor do they persist once the function has run. These are called “local” variables, and are only accessible inside their function. However, “global” variables are visible inside and outside of functions. In python, you can assign global variables. Type the following script in `scope.py` and try it:

```
## Try this first

_a_global = 10

def a_function():
    _a_global = 5
    _a_local = 4
    print "Inside the function, the value is ", _a_global
    print "Inside the function, the value is ", _a_local
    return None

a_function()
print "Outside the function, the value is ", _a_global


## Now try this

_a_global = 10

def a_function():
    global _a_global
    _a_global = 5
    _a_local = 4
    print "Inside the function, the value is ", _a_global
    print "Inside the function, the value is ", _a_local
    return None

a_function()
print "Outside the function, the value is", _a_global
```

However, in general, avoid assigning globals because you run the risk of “exposing” unwanted variables to all functions within your name work/namespace.

### 5.8.4 Running python code from terminal or bash shell

We have been running scripts from our beloved terminal or bash shell. To execute script file from within the python shell, use `execfile("filetoload.py")`. In IPython, use: `run filetoload.py`

## 5.9 Control statements

OK, let's get deeper into python functions. To begin, first copy and rename `boilerplate.py` (to make use of it's existing structure and save you some typing):

```
$ cp boilerplate.py control_flow.py
```

Then type the following script into `control_flow.py`:

```
#!/usr/bin/python

"""Some functions exemplifying the use of control statements"""
#docstrings are considered part of the running code (normal comments are
#stripped). Hence, you can access your docstrings at run time.
__author__ = 'Samraat Pawar (s.pawar@imperial.ac.uk)'
__version__ = '0.0.1'

import sys

def even_or_odd(x=0): # if not specified, x should take value 0.

    """Find whether a number x is even or odd."""
    if x % 2 == 0: #The conditional if
        return "%d is Even!" % x
    return "%d is Odd!" % x

def largest_divisor_five(x=120):
    """Find which is the largest divisor of x among 2,3,4,5."""
    largest = 0
    if x % 5 == 0:
        largest = 5
    elif x % 4 == 0: #means "else, if"
        largest = 4
    elif x % 3 == 0:
        largest = 3
    elif x % 2 == 0:
        largest = 2
    else: # When all other (if, elif) conditions are not met
        return "No divisor found for %d!" % x # Each function can return a ←
            value or a variable.
    return "The largest divisor of %d is %d" % (x, largest)

def is_prime(x=70):
    """Find whether an integer is prime."""
    for i in range(2, x): # "range" returns a sequence of integers
        if x % i == 0:
            print "%d is not a prime: %d is a divisor" % (x, i) #Print ←
                formatted text "%d %s %f %e" % (20,"30",0.0003,0.00003)
```

```

        return False
print "%d is a prime!" % x
return True

def find_all_primes(x=22):
    """Find all the primes up to x"""
    allprimes = []
    for i in range(2, x + 1):
        if is_prime(i):
            allprimes.append(i)
    print "There are %d primes between 2 and %d" % (len(allprimes), x)
    return allprimes

def main(argv):
    # sys.exit("don't want to do this right now!")
    print even_or_odd(22)
    print even_or_odd(33)
    print largest_divisor_five(120)
    print largest_divisor_five(121)
    print is_prime(60)
    print is_prime(59)
    print find_all_primes(100)
    return 0

if (__name__ == "__main__"):
    status = main(argv)
    sys.exit(status)

```

Note the “docstrings” are in the above code enclosed in “““ “““. A docstring is meant to describe the operation of the function or script. Docstrings are considered part of the running code, while normal comments are stripped. Hence, you can access your docstrings at run time. It is a good idea to have docstrings at the start of every python script and module as it can provide useful information to the user and you as well, down the line. Now run the code:

```
In []: run control_flow.py
```

You can also call any of the functions within `control_flow.py`:

```
In []: even_or_odd(11)
Out[]: '11 is Odd!'
```

This is possible without explicitly importing the modules because you are only running one script. You would have to do an explicit `import` if you needed a module from another python script file.

### 5.9.1 Control flow exercises

Write the following, and save them to `cfexercises.py`.

```
# How many times will 'hello' be printed?
# 1)
for i in range(3, 17):
    print 'hello'
```

```

# 2)
for j in range(12):
    if j % 3 == 0:
        print 'hello'

# 3)
for j in range(15):
    if j % 5 == 3:
        print 'hello'
    elif j % 4 == 3:
        print 'hello'

# 4)
z = 0
while z != 15:
    print 'hello'
    z = z + 3

# 5)
z = 12
while z < 100:
    if z == 31:
        for k in range(7):
            print 'hello'
    elif z == 18:
        print 'hello'
    z = z + 1

# What does fooXX do?
def foo1(x):
    return x ** 0.5

def foo2(x, y):
    if x > y:
        return x
    return y

def foo3(x, y, z):
    if x > y:
        tmp = y
        y = x
        x = tmp
    if y > z:
        tmp = z
        z = y
        y = tmp
    return [x, y, z]

def foo4(x):
    result = 1
    for i in range(1, x + 1):
        result = result * i
    return result

# This is a recursive function, meaning that the function calls itself
# read about it at
# en.wikipedia.org/wiki/Recursion_(computer_science)
def foo5(x):
    if x == 1:
        return 1

```

```

    return x * foo5(x - 1)
foo5(10)

```

## 5.10 Loops

Write the following, and save them to loops.py.

```

# for loops in Python
for i in range(5):
    print i

my_list = [0, 2, "geronimo!", 3.0, True, False]
for k in my_list:
    print k

total = 0
summands = [0, 1, 11, 111, 1111]
for s in summands:
    print total + s

# while loops in Python
z = 0
while z < 100:
    z = z + 1
    print (z)

b = True
while b:
    print "GERONIMO! infinite loop! ctrl+c to stop!"
    # ctrl + c to stop!

```

### 5.10.1 List comprehensions

Python offers a way to combine loops, functions and logical tests in a single line of code. Type the following in a script file called oaks.py:

```

## Let's find just those taxa that are oak trees from a list of species

taxa = [ 'Quercus robur',
          'Fraxinus excelsior',
          'Pinus sylvestris',
          'Quercus cerris',
          'Quercus petraea',
      ]

def is_an_oak(name):
    return name.lower().startswith('quercus ')

##Using for loops
oaks_loops = set()
for species in taxa:
    if is_an_oak(species):

```

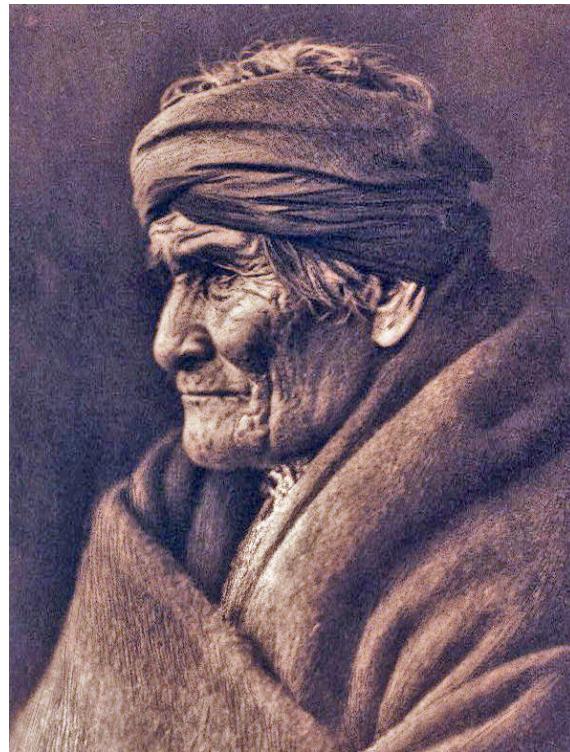


Figure 5.3: In case you were wondering who Geronimo was.

```

        oaks_loops.add(species)
print oaks_loops

##Using list comprehensions
oaks_lc = set([species for species in taxa if is_an_oak(species)])
print oaks_lc

##Get names in UPPER CASE using for loops
oaks_loops = set()
for species in taxa:
    if is_an_oak(species):
        oaks_loops.add(species.upper())
print oaks_loops

##Get names in UPPER CASE using list comprehensions
oaks_lc = set([species.upper() for species in taxa if is_an_oak(species)])
print oaks_lc

```

Don't go mad with list comprehensions — code readability is more important than squeezing lots into a single line!

## 5.11 Practical

1. Review and make sure you can run all the commands, code fragments, and scripts we have covered today and get the expected outputs.
2. Run `boilerplate.py` and `control_flow.py` from the bash terminal instead of from

within the ipython shell (try both python and ipython from the bash)

3. Take `cfexercises.py` and convert it into a function such that all the `fooXX` functions can take arguments from the user (like the modules/functions inside `control_flow.py`). Also Add some test arguments to show that they work (again, like `control_flow.py`) — for example, “`foo5(10)`”. Thus, running `cfexercises.py` should now also output evaluations of all the `fooXX` modules along with a bunch of hellos.
4. Open and complete the tasks in `lc1.py`, `lc2.py`, `dictionary.py`, `tuple.py` (you can tackle them in any order)
5. Include an appropriate docstring (if one is missing) at the beginning of *each* of each of the python script / module files you have written, as well as at the start of every function (or sub-module) in a module.
6. Also annotate your code lines as much and as often as necessary using `#`.
7. Keep all code files organized in `CMEECourseWork/Week2/Code` and git push all your work. The updated bitbucket git repository should contain: `basic_csv.py`, `basic_io.py`, `boilerplate.py`, `sysargv.py`, `using_name.py`, `control_flow.py`, `scope.py`, `loops.py`, `cfexercises.py`, `oaks.py`, `lc1.py`, `lc2.py`, `dictionary.py`, & `tuple.py` — all scripts should work on any other linux laptop!

## 5.12 Readings and Resources

- The Zen of python: open a python shell and type `import this`
- Code like a Pythonista: Idiomatic python (Google it)
- Also good: the Google python Style Guide
- Browse the python tutorial: [www.docs.python.org/tutorial/](http://www.docs.python.org/tutorial/)
- For functions and modules:  
[www.learnpythononthehardway.org/book/ex40.html](http://www.learnpythononthehardway.org/book/ex40.html)
- For IPython:  
[www.ipython.org/ipython-doc/stable/interactive/tips.html](http://www.ipython.org/ipython-doc/stable/interactive/tips.html) [www.wiki.ipython.org/Cookbook](http://www.wiki.ipython.org/Cookbook)



# Chapter 6

## Biological Computing in Python II

...some things in life are bad. They can really make you mad. Other things just make you swear and curse. When you're chewing on life's gristle, don't grumble; give a whistle, and this'll help things turn out for the best. And... always look on the bright side of life...

---

—*Guess who?*

### 6.1 ipython

ipython stands for interactive python shell. In IPython, you can TAB everything — a lot will be revealed!

In particular, TAB leads to autocompletion. Also TAB after “.” reveals object’s functions and attributes:

```
In []: alist = ['a', 'b', 'c']

In []: alist. #TAB now!
alist.append    alist.extend    alist.insert    alist.remove
alist.sort      alist.count     alist.index     alist.pop
alist.reverse

In []: adict = {'mickey': 100, 'mouse': 3.14}

In []: adict.
adict.clear      adict.items      adict.pop
adict.viewitems  adict.copy       adict.iteritems
adict.popitem    adict.viewkeys   adict.fromkeys
adict.iterkeys   adict.setdefault adict.viewvalues
adict.get        adict.itervalues adict.update
adict.has_key    adict.keys      adict.values
```

IPython also has “magic commands” (start with %; e.g., %run). Some useful magic commands:

%who	Shows current namespace (all variables, modules and functions)
%whos	Also display the type of each variable; typing %whos function only displays functions etc.
%pwd	Print working directory
%history	Print recent commands
%cpaste	Paste indented code into IPython — very useful when you want to run just part of a function (indentation and all)

Let's try the `%cpaste` function. First, type the following code in a temporary file:

```
def PrintNumbers(x):  
    for i in range(x):  
        print x  
    return 0
```

Now launch `ipython` and type:

```
In []: x = 11

# Now copy the for loop from the file.
# Note that there is extra indentation!

In []: %cpaste
Pasting code; enter '---' alone on the line to stop
or use Ctrl-D.
:     for i in range(x):
:         print x
:--
11
11
11
11
11
11
11
11
11
11
11
11
11
11
11
```

Another useful feature is the question mark, which can be used to find what a particular Python object is, including variables you created! For example, try this to check what `adict`, which you created above, is:

```
In []: ?adict
Type:      dict
Base Class: <type 'dict'>
String Form:{'mickey': 100, 'mouse': 3.14}
Namespace:  Interactive
Length:     2
Docstring:
dict() -> new empty dictionary
dict(mapping) -> new dictionary initialized from a mapping
object's
    (key, value) pairs
```

```

dict(iterable) -> new dictionary initialized as if via:
    d = {}
    for k, v in iterable:
        d[k] = v
dict(**kwargs) -> new dictionary initialized with the
name=value pairs
    in the keyword argument list.
For example: dict(one=1, two=2)

```

## 6.2 Functions, Modules, and code compartmentalization

Ideally you should aim to compartmentalize your code into a bunch of functions, typically written in a single .py file: this are Python “modules”, which you were introduced to previously. Why bother with modules? Because:

- Keeping code compartmentalized is good for debugging, unit testing, and profiling (coming up later)
- Makes code more compact by minimizing redundancies (write repeatedly used code segments as a module)
- Allows you to import and use useful functions that you yourself wrote, just like you would from standard python packages (see next slide)

### 6.2.1 Importing Modules

There are different ways to **import** a module:

- `import my_module`, then functions in the module can be called as `my_module.one_of_my_functions()`.
- `from my_module import my_function` imports only the function `my_function` in the module `my_module`. It can then be called as if it were part of the main file: `my_function()`.
- `import my_module as mm` imports the module `my_module` and calls it `mm`. Convenient when the name of the module is very long. The functions in the module can be called as `mm.one_of_my_functions()`.
- `from my_module import *`. Avoid doing this!  
*Why?* – to avoid name conflicts!
- You can also access variables written into modules: `import my_module, then my_module.one_of_my_variables`

## 6.3 Python packages

A Python package is simply a directory of Python modules (quite like an R package). Many packages, such as the following that I find particularly useful, are always available as standard libraries (just require `import` from within python or ipython):

- `io`: file input-output with `*.csv`, `*.txt`, etc.
- `subprocess`: to run other programs, including multiple ones at the same time, including operating system-dependent functionality

- `sqlite3`: for manipulating and querying `sqlite` databases
- `math`: for mathematical functions

A million (literally?) other are accessible by explicitly installing them using  
`sudo apt-get install python-packagename` (as you did for `ipdb` and `scipy` previously). Some particularly mentionable ones are:

- `scipy`: for scientific computing
- `matplotlib`: for plotting (very matlab-like, requires `scipy`) (all packaged in `pylab`)
- `scrapy`: for writing web spiders that crawl web sites and extract data from them
- `beautifulsoup`: for parsing HTML and XML (can do what `scrapy` does)
- `biopython`: for bioinformatics

Of course, you have already installed some of these (`scipy`, `matplotlib`).

For those of you interested in bioinformatics, the `biopython` package will be particularly useful. We will not cover bioinformatics in any depth within the python weeks, but you will have the (mostly implicit) opportunity to use Python for bioinformatics in other weeks, especially the Genomics weeks, and perhaps use it for your Masters project. Therefore, I suggest that if bioinformatics is your thing, check out `biopython` — in particular the worked examples at <http://biopython.org/DIST/docs/tutorial/Tutorial.html>.

Let's shift gears now, and look at a very important skill that you should learn, or at least be aware of — *Regular expressions*.

## 6.4 Regular expressions in python

Regular expressions (regex) are a tool to find patterns in strings, such as:

- Find DNA motifs in sequence data
- Navigate through files in a directory
- Parse text files
- Extract information from html and xml files

Regex packages are available for most programming languages (`grep` in UNIX / Linux, where regex first became popular).

A regex may consist of a combination of “metacharacters” (modifiers) and “regular” or literal characters. There are 14 metacharacters: `[ ] { } () \ ^ $ . | ? * +`

These metacharacters do special things, for example:

- `[12]` means match target to 1 and if that does not match then match target to 2
- `[0-9]` means match to any character in range 0 to 9
- `[^Ff]` means anything except upper or lower case f and `[^a-z]` means everything except lower case a to z

Everything else is interpreted literally (e.g., `a` is matched by entering `a` in the regex). A useful (not exhaustive) list of regex elements is:

a	match the character a
3	match the number 3
\n	match a newline
\t	match a tab
\s	match a whitespace
.	match any character except line break (newline)
\w	match any alphanumeric character (including underscore)
\W	match any character not covered by \w (i.e., match any non-alphanumeric character excl. underscore)
\d	match a numeric character
\D	match any character not covered by \d (i.e., match a non-digit)
[atgc]	match any character listed: a, t, g, c
at gc	match at or gc
[^atgc]	any character not listed: any character but a, t, g, c
?	match the preceding pattern element zero or one times
*	match the preceding pattern element zero or more times
+	match the preceding pattern element one or more times
{n}	match the preceding pattern element exactly n times
{n,}	match the preceding pattern element at least n times
{n,m}	match the preceding pattern element at least n but not more than m times
^	match the beginning of a line
\$	match the end of a line

Regex functions in python are in the module `re` — so we will import `re`. The simplest python regex function is `re.search`, which searches the string for match to a given pattern — returns a *match object* if a match is found and `None` if not.

**Always** put `r` in front of your regex — it tells python to read the regex in its “raw” (literal) form. Without raw string notation (`r“text”`), every backslash ('\') in a regular expression would have to be prefixed with another one to escape it.

From <https://docs.python.org/2/library/re.html>: *If you’re not using a raw string to express the pattern, remember that Python also uses the backslash as an escape sequence in string literals; if the escape sequence isn’t recognized by Python’s parser, the backslash and subsequent character are included in the resulting string. However, if Python would recognize the resulting sequence, the backslash should be repeated twice. This is complicated and hard to understand, so it’s highly recommended that you use raw strings for all but the simplest expressions.*

OK, let’s try some regexes (type all that follows in `Code/regexes.py`):

```
import re

my_string = "a given string"
# find a space in the string
match = re.search(r'\s', my_string)

print match
# this should print something like
# <_sre.SRE_Match object at 0x93ecdd30>

# now we can see what has matched
```

```

match.group()

match = re.search(r's\w*', my_string)

# this should return "string"
match.group()

# NOW AN EXAMPLE OF NO MATCH:
# find a digit in the string
match = re.search(r'\d', my_string)

# this should print "None"
print match

# Further Example
#
my_string = 'an example'
match = re.search(r'\w*\s', my_string)

if match:
    print 'found a match:', match.group()
else:
    print 'did not find a match'

```

To know whether a pattern was matched, we can use an `if`:

```

MyStr = 'an example'

match = re.search(r'\w*\s', MyStr)

if match:
    print 'found a match:', match.group()
else:
    print 'did not find a match'

```

Here are some more regexes (add all that follows to the `Code/regexes.py`):

```

# Some Basic Examples
match = re.search(r'\d', "it takes 2 to tango")
print match.group() # print 2

match = re.search(r'\s\w*\s', 'once upon a time')
match.group() # ' upon '

match = re.search(r'\s\w{1,3}\s', 'once upon a time')
match.group() # ' a '

match = re.search(r'\s\w*$', 'once upon a time')
match.group() # ' time'

match = re.search(r'\w*\s\d.*\d', 'take 2 grams of H2O')
match.group() # 'take 2 grams of H2'

match = re.search(r'^\w*.*\s', 'once upon a time')
match.group() # 'once upon a '
## NOTE THAT *, +, and { } are all "greedy":
## They repeat the previous regex token as many times as possible
## As a result, they may match more text than you want

```

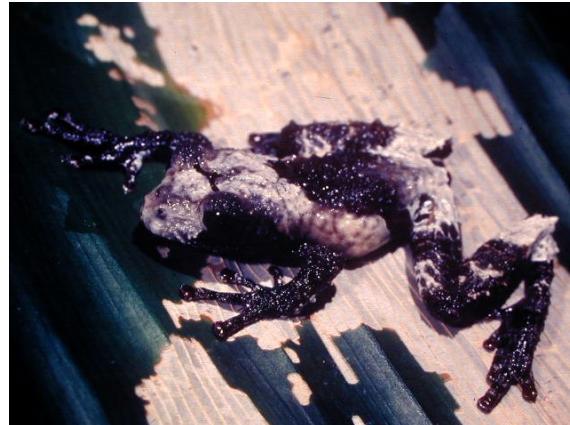


Figure 6.1: In case you were wondering what *Theloderma asper*, the bird-shit frog, looks like. I snapped this one in North-east India ages ago

```
## To make it non-greedy, use ?:  
match = re.search(r'^\w*.*?\s', 'once upon a time')  
match.group() # 'once '  
  
## To further illustrate greediness, let's try matching an HTML tag:  
match = re.search(r'<.+>', 'This is a <EM>first</EM> test')  
match.group() # '<EM>first</EM>'  
## But we didn't want this: we wanted just <EM>  
## It's because + is greedy!  
  
## Instead, we can make + "lazy"!  
match = re.search(r'<.+?>', 'This is a <EM>first</EM> test')  
match.group() # '<EM>'  
  
## OK, moving on from greed and laziness  
match = re.search(r'\d*\.\?\d*', '1432.75+60.22i') #note "\" before "."  
match.group() # '1432.75'  
  
match = re.search(r'\d*\.\?\d*', '1432+60.22i')  
match.group() # '1432'  
  
match = re.search(r'[AGTC]+', 'the sequence ATTCGT')  
match.group() # 'ATTCGT'  
  
re.search(r'\s+[A-Z]\{1\}\w+\s\w+', 'The bird-shit frog''s name is Theloderma ↫  
asper').group() # ' Theloderma asper'  
## NOTE THAT I DIRECTLY RETURNED THE RESULT BY APPENDING .group()
```

You can group regexes into meaningful blocks using parentheses. For example, let's try matching a string consisting of an academic's name, email address and research area or interest (no need to type this into any python file):

```
MyStr = 'Samraat Pawar, s.pawar@imperial.ac.uk, Systems biology and  
ecological theory'  
  
# without groups  
match = re.search(r"[\w\s]*,\s[\w\.\@]*,\s[\w\s\&]*", MyStr)
```

```

match.group()
'Samraat Pawar, s.pawar@imperial.ac.uk, Systems biology and ecological theory←
'

match.group(0)
'Samraat Pawar, s.pawar@imperial.ac.uk, Systems biology and ecological theory←
'

# now add groups using ( )
match = re.search(r"([\w\s]*),\s([\w\.\@]*),\s([\w\s&]*)",MyStr)

match.group(0)
'Samraat Pawar, s.pawar@imperial.ac.uk, Systems biology and ecological theory←
'

match.group(1)
'Samraat Pawar'

match.group(2)
's.pawar@imperial.ac.uk'

match.group(3)
'Systems biology and ecological theory'

```

Have a look at `re4.py` in your code repository for more in parsing email addresses using regexes.

#### 6.4.1 Some RegExercises

These exercises are not for submission as part of your coursework, but we will discuss them in class.

1. Translate the following regular expressions into regular English (don't type this in `regexprs.py`)!

```

r'^abc[ab]+\\s\\t\\d'
r'^\\d{1,2}\\/\\d{1,2}\\/\\d{4}$'
r'\\s*[a-zA-Z,\\s]+\\s*'

```

2. Write a regex to match dates in format YYYYMMDD, making sure that:

- Only seemingly valid dates match (i.e., year greater than 1900)
- First digit in month is either 0 or 1
- First digit in day  $\leq 3$



[www.xkcd.com/208/](http://www.xkcd.com/208/)

#### 6.4.2 Important re functions

`re.compile(reg)`

Compile a regular expression. In this way the pattern is stored for repeated use, improving the speed.

`re.search(reg, text)`

Scan the string and find the first match of the pattern in the string. Returns a `match` object if successful and `None` otherwise.

`re.match(reg, text)`

as `re.search`, but only match the beginning of the string.

`re.split(ref, text)`

Split the text by the occurrence of the pattern described by the regular expression.

`re.findall(ref, text)`

As `re.search`, but return a list of all the matches. If groups are present, return a list of groups.

`re.finditer(ref, text)`

As `re.search`, but return an iterator containing the next match.

`re.sub(ref, repl, text)`

Substitute each non-overlapping occurrence of the match with the text in `repl` (or a function!).

## 6.5 Numerical computing in python

The python package `scipy` can help you do serious number crunching including,

- Linear algebra
- Numerical integration
- Fourier transforms
- Interpolation
- Special functions (Incomplete Gamma, Bessel, etc.)

- Generation of random numbers
- Statistical functions

Let's look at `scipy`. In the following, we will use the `array` class in `scipy` for data manipulations and calculations. Scipy arrays are similar in some respects to python lists, but are more naturally multidimensional, homogeneous in type, and allow efficient manipulations. So try this:

```
In [1]: import scipy

In [2]: a = scipy.array(range(5))

In [3]: a
Out[3]: array([0, 1, 2, 3, 4])

In [4]: a = scipy.array(range(5), dtype = float)

In [5]: a
Out[5]: array([ 0.,  1.,  2.,  3.,  4.])

In [6]: a.dtype
Out[6]: dtype('float64')

In [7]: x = scipy.arange(5)

In [8]: x
Out[8]: array([0, 1, 2, 3, 4])

In [9]: x = scipy.arange(5.)

In [10]: x
Out[10]: array([ 0.,  1.,  2.,  3.,  4.])

In [11]: x. # TAB after x. to see methods you can apply to x
x.T          x.conj         x.fill
x nbytes     x.round        x.take
x.all        x.conjugate    x.flags
x.ndim       x.searchsorted x.tofile
x.any        x.copy         x.flat
x.newbyteorder x.setfield   x.tolist
x.argmax     x.ctypes      x.flatten
x.nonzero    x.setflags    x.tostring
x.argmin     x.cumprod    x.getfield
x.prod       x.shape       x.trace
x.argsort    x.cumsum     x.imag
x.ptp        x.size        x.transpose
x.astype     x.data        x.item
x.put        x.sort        x.var
x.base       x.diagonal   x.itemset
x.ravel      x.squeeze    x.view
x.byteswap   x.dot         x.itemsize
x.real       x.std         x.max
x.choose     x.dtype       x.mean
x.repeat     x.strides
x.clip       x.dump
x.reshape    x.sum
x.compress  x.dumps
x.resize    x.swapaxes
```

```
In [11]: x.tolist()
```

```
Out[11]: [0.0, 1.0, 2.0, 3.0, 4.0]

In [12]: x.shape
Out[12]: (5,)

In [14]: mat = scipy.array([[0, 1], [2, 3]])

In [16]: mat.shape
Out[16]: (2, 2)

In [17]: mat[1]
Out[17]: array([2, 3])

In [18]: mat[0,0]
Out[18]: 0

In [19]: mat.ravel() # flatten!
Out[19]: array([0, 1, 2, 3])

In [20]: scipy.ones((4,2))
Out[20]:
array([[ 1.,  1.],
       [ 1.,  1.],
       [ 1.,  1.],
       [ 1.,  1.]))

In [21]: scipy.zeros((4,2))
Out[21]:
array([[ 0.,  0.],
       [ 0.,  0.],
       [ 0.,  0.],
       [ 0.,  0.]]))

In [22]: scipy.identity(4)
Out[22]:
array([[ 1.,  0.,  0.,  0.],
       [ 0.,  1.,  0.,  0.],
       [ 0.,  0.,  1.,  0.],
       [ 0.,  0.,  0.,  1.]]))

In [23]: m = scipy.identity(4)

In [24]: m.reshape((8, 2))
Out[24]:
array([[ 1.,  0.],
       [ 0.,  0.],
       [ 0.,  1.],
       [ 0.,  0.],
       [ 0.,  0.],
       [ 1.,  0.],
       [ 0.,  0.],
       [ 0.,  1.]))

In [25]: m.fill(16)

In [26]: m
Out[26]:
array([[ 16.,  16.,  16.,  16.],
       [ 16.,  16.,  16.,  16.],
       [ 16.,  16.,  16.,  16.],
       [ 16.,  16.,  16.,  16.]])
```

Let's perform some common operations on arrays:

```
In [1]: import scipy

In [2]: mm = scipy.arange(16)

In [3]: mm = mm.reshape(4,4)

In [4]: mm
Out[4]:
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11],
       [12, 13, 14, 15]])

# NOTE: Rows first

In [5]: mm.transpose()
Out[5]:
array([[ 0,  4,  8, 12],
       [ 1,  5,  9, 13],
       [ 2,  6, 10, 14],
       [ 3,  7, 11, 15]])

In [6]: mm + mm.transpose()
Out[6]:
array([[ 0,  5, 10, 15],
       [ 5, 10, 15, 20],
       [10, 15, 20, 25],
       [15, 20, 25, 30]])

In [7]: mm - mm.transpose()
Out[7]:
array([[ 0, -3, -6, -9],
       [ 3,  0, -3, -6],
       [ 6,  3,  0, -3],
       [ 9,  6,  3,  0]])

In [8]: mm * mm.transpose()
## Elementwise!

Out[8]:
array([[ 0,   4,  16,  36],
       [ 4,  25,  54,  91],
       [16,  54, 100, 154],
       [36,  91, 154, 225]])

In [9]: mm / mm.transpose()
Warning: divide by zero encountered in divide

# Note the integer division
Out[9]:
array([[0, 0, 0, 0],
       [4, 1, 0, 0],
       [4, 1, 1, 0],
       [4, 1, 1, 1]])

In [10]: mm * scipy.pi
Out[10]:
array([[ 0.          ,    3.14159,    6.28318531,   9.42477796],
```

```
[ 12.566370,  15.70796,  18.84955592,  21.99114858],
[ 25.132741,  28.27433,  31.41592654,  34.55751919],
[ 37.699111,  40.84070,  43.98229715,  47.1238898 ]])

In [11]: mm.dot(mm) # MATRIX MULTIPLICATION
Out[11]:
array([[ 56,   62,   68,   74],
       [152,  174,  196,  218],
       [248,  286,  324,  362],
       [344,  398,  452,  506]])
```

We can do a lot more (but won't!) by importing the `linalg` sub-package: `scipy.linalg`.

Two other particularly useful `scipy` sub-packages are `scipy.integrate` (what will I need this for?) and `scipy.stats` (why not use R for this?). Let's take a quick spin in `scipy.stats`!

```
In [18]: import scipy.stats

In [19]: scipy.stats.
scipy.stats.arcsine          scipy.stats.lognorm
scipy.stats.bernoulli        scipy.stats.mannwhitneyu
scipy.stats.beta             scipy.stats.maxwell
scipy.stats.binom            scipy.stats.moment
scipy.stats.chi2              scipy.stats.nanstd
scipy.stats.chisqprob        scipy.stats.nbinom
scipy.stats.circvar          scipy.stats.norm
scipy.stats.expon             scipy.stats.powerlaw
scipy.stats.gompertz          scipy.stats.t
scipy.stats.kruskal           scipy.stats.uniform

In [19]: scipy.stats.norm.rvs(size = 10) # 10 samples from
N(0,1)
Out[19]:
array([-0.951319, -1.997693,  1.518519, -0.975607,  0.8903,
       -0.171347, -0.964987, -0.192849,  1.303369,  0.6728])

In [20]: scipy.stats.norm.rvs(5, size = 10)
# change mean to 5
Out[20]:
array([ 6.079362,  4.736106,  3.127175,  5.620740,  5.98831,
       6.657388,  5.899766,  5.754475,  5.353463,  3.24320])

In [21]: scipy.stats.norm.rvs(5, 100, size = 10)
# change sd to 100
Out[21]:
array([-57.886247,   12.620516,  104.654729,  -30.979751,
       41.775710,  -31.423377,  -31.003134,   80.537090,
       3.835486,   103.462095])

# Random integers between 0 and 10
In [23]: scipy.stats.randint.rvs(0, 10, size =7)
Out[23]: array([6, 6, 2, 0, 9, 8, 5])
```

### 6.5.1 A numerical integration example

Let's look at an example using `scipy.integrate`. Create `LV1.py` in your Week2/Code directory and run it. The Lotka-Volterra model is:

$$\frac{dR}{dt} = rR - aCR \frac{dC}{dt} = -zC + eaCR \quad (6.1)$$

where  $C$  and  $R$  are consumer (e.g., predator) and resource (e.g., prey) population sizes (either biomass or numbers),  $r$  is the intrinsic growth rate of the resource population,  $a$  is a “search rate” that determines the encounter rate between consumer and resource,  $z$  is mortality rate and  $e$  is the consumer’s efficiency in converting resource to consumer biomass.

`LV1.py` runs (numerically solves the ODE) this model and plots the equilibrium. Have good look at the code, line by line, and make sure that you understand what’s going on. A subsequent practical will require you to use this code to extend the LV model.

## 6.6 Practical

### 6.6.1 Align DNA sequences

Align two DNA sequences such that they are as similar as possible.

The idea is to start with the longest string and try to position the shorter string in all possible positions. For each position, count a “score” : number of bases matched perfectly over the number of bases attempted. Your tasks:

1. Open and run `Practicals/Code/align_seqs.py` and make sure you understand what each line is doing (hint: you can use `pdb` to do this)
2. Convert `align_seqs.py` to a Python function that takes the DNA sequences as an input from a `.csv` file and saves the best alignment along with its corresponding score in a single text file (hint: remember `pickle`!)
3. Make sure you provide a test `.csv` file that the script can call
4. Make sure you provide a test `.csv` file that the script can call.
5. Don’t forget to add one or more doctrings
6. Extra Credit – align the `.fasta` sequences from Week 1; call the new script `align_seqs.fasta.py`.

### 6.6.2 Blackbirds problem

Complete the code `blackbirds.py` that you find in the CMEEMasteRepo (necessary data file is also there).

## 6.7 Readings and Resources

- [docs.python.org/2/howto/regex.html](https://docs.python.org/2/howto/regex.html)

- Google's short class on regex in python:  
<https://developers.google.com/edu/regular-expressions>
- [www.regular-expressions.info](http://www.regular-expressions.info) has a good intro, tips and a great array of canned solutions
- Use and abuse of regex:  
[www.codinghorror.com/blog/2005/02/regex-use-vs-regex-abuse.html](http://www.codinghorror.com/blog/2005/02/regex-use-vs-regex-abuse.html)
- [www.matplotlib.org/](http://www.matplotlib.org/)
- For SciPy, the official documentation is best:  
[docs.scipy.org/doc/scipy/reference/](http://docs.scipy.org/doc/scipy/reference/)  
Read about the scipy modules you think will be important to you...
- Many illustrative examples at <http://wiki.scipy.org/Cookbook> (including Lotka-Volterra!)
- In general, good module-specific cookbooks are out there (e.g., biopython)



# Chapter 7

## Additional Python topics, and a Long practical

### 7.1 Using python to build workflows

You can use python to build an automated data analysis or simulation workflow that involves multiple applications, especially the ones you have already learnt: R, L<sup>A</sup>T<sub>E</sub>X, & UNIX bash. For example, you could, in theory, write a single Python script to generate and update your masters dissertation, tables, plots, and all. Python is ideal for building such workflows because it has packages for practically every purpose (see Section on Packages above).

### 7.2 Using subprocess

The `subprocess` module is particularly important as it can run other applications, including R. Let's try – first launch `ipython`, then `cd` to your python code directory, and type:

```
import subprocess
subprocess.os.system("geany boilerplate.py")
subprocess.os.system("gedit ../Data/TestOaksData.csv")
subprocess.os.system("python boilerplate.py") # A bit silly!
```

Easy as pie! Similarly, to compile your Latex document (using `pdflatex` in this case):

```
subprocess.os.system("pdflatex yourlatexdoc.tex")
```

You can also do this (instead of using `subprocess.os`):

```
subprocess.Popen("geany boilerplate.py", shell=True).wait()
```

You can also use `subprocess.os` to make your code OS (Linux, Windows, Mac) independent. For example to assign paths:

```
subprocess.os.path.join('directory', 'subdirectory', 'file')
```

The result would be appropriately different on Windows (with backslashes instead of forward slashes).

Note that in all cases you can “catch” the output of subprocess so that you can then use the output within your python script. A simple example, where the output is a platform-dependent directory path, is:

```
MyPath = subprocess.os.path.join('directory', 'subdirectory', 'file')
```

Explore what subprocess can do by tabbing subprocess., and also for submodules, e.g., type subprocess.os. and then tab.

### 7.2.1 Running R

R is likely an important part of your workflow, for example for statistical analyses and pretty plotting (hmmm... ggplot2!). Try the following.

Create an R script file called TestR.R in your CMEECourseWork/Week6/Code with the following content:

```
print("Hello, this is R!")
```

Then, create TestR.py in CMEECourseWork/Week6/Code with the following content :

```
import subprocess
subprocess.Popen("/usr/lib/R/bin/Rscript --verbose TestR.R > \
.../Results/TestR.Rout 2> .../Results/TestR_errFile.Rout", \
shell=True).wait()
```

*Note the backslashes — this is so that python can read the multiline script as a single line.*

It is possible that the location of RScript is different in your Ubuntu install. To locate it, try find /usr -name 'Rscript' in the linux terminal (not in python!).

Now run TestR.py (or %cpaste) and check TestR.Rout and TestR\_errorFile.Rout.

Also check what happens if you run (type directly in ipython or python console):

```
subprocess.Popen("/usr/lib/R/bin/Rscript --verbose NonExistScript.R > \
.../Results/outputFile.Rout 2> .../Results/errorFile.Rout", \
shell=True).wait()
```

What do you see on the screen? Now check outputFile.Rout and errorFile.Rout.

## 7.3 Long Practical – Building a workflow

We have talked a lot about workflows and confronting models with data. It's time to do something concrete with all the stuff you have learnt!

This practical gives you an opportunity to try the “whole nine yards” of developing and implementing a workflow and delivering a “finished product” to answer a scientific question in biology.

**This practical is optional for the MRes’, but can provide extra credit if you choose to participate!**

*Huh? how will this work?* — One MSc and *up to two* MRes’ can collaborate and co-author the report. The MSc is expected to do the bulk of the work, of course! The contributions will be stated in the report (see report subsection below). Every participating MRes, concomitant with her/his contribution, gets some extra credit that can go towards bolstering weaknesses in the weekly practical assessments. I cannot give any more details, because I have never tried this before!

### 7.3.1 Background and Objectives

**The main question you will address is:** *How well does a mechanistic model based upon biochemical principles fit a dataset of thermal responses of individual fitness in phytoplankton?*

This is currently a “hot” (no pun intended) topic in biology, with both ecological and evolutionary consequences, as we discussed in the earlier lecture. On the *ecological side*, because the temperature-dependence of metabolic rate sets the rate of intrinsic  $r_{max}$  (papers by Savage et al., Brown et al.) as well as interactions between species, it has a strong effect on population dynamics. In this context, note that 99.9% of life on earth is ectothermic! On the *evolutionary side*, the temperature-dependence of fitness and species interactions also means that warmer environments may have stronger rates of evolution. This may be compounded by the fact that mutation rates may also increase with temperature (papers by Gillooly et al.)!

### 7.3.2 The Data

These data are for a subset of growth rates ( $\text{time}^{-1}$ ) of phytoplankton, collected by experimental studies from across the world, and compiled by Thomas et al. and Chen et al. The figure in your `Data` folder shows an example curve along with a fitted curve of the Schoolfield model.

### 7.3.3 The Models

There are at least two alternative models.

The Schoolfield model (paper is in `Readings` directory) is a mechanistic alternative:

$$B(T) = B_0 \cdot \frac{e^{\left[ \frac{-E}{k} \left( \frac{1}{T} - \frac{1}{283.15} \right) \right]}}{1 + \frac{E}{E_D - E} \cdot e^{\left[ \frac{E_D}{k} \left( \frac{1}{T_{pk}} - \frac{1}{T} \right) \right]}} \quad (7.1)$$

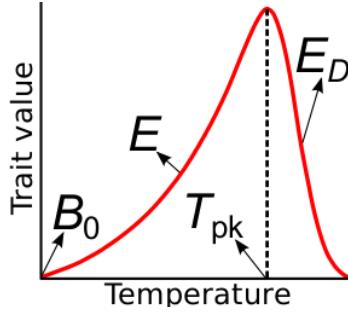


Figure 7.1: Illustration of a unimodal thermal response curve for a biological trait according to the Schoolfield model (Equation 7.1).  $B_0$  is the trait value at low temperature and controls the offset of the curve.  $E$  is the activation energy (eV) which controls the rise of the curve up to the peak, whereas  $E_D$  is the de-activation energy (eV) which controls the fall of the curve after the peak.  $T_{pk}$  (K) is the temperature where trait performance peaks ( $B_{pk}$ ).

Here,  $k$  is the Boltzmann constant ( $8.617 \times 10^{-5} \text{ eV} \cdot \text{K}^{-1}$ ),  $B$  the value of the trait at a given temperature  $T$  (K) ( $K = {}^\circ\text{C} + 273.15$ ), while  $B_0$  is the trait value at 283.15 K ( $10^\circ\text{C}$ ) which stands for the value of the growth rate at low temperature and controls the vertical offset of the curve.  $E$  is the activation energy (eV) which controls the rise of the curve up to the peak,  $E_D$  is the de-activation energy (eV) which controls the fall of the curve after the peak and  $T_{pk}$  (K) is the temperature where trait performance is maximal.

The Gaussian-Gompertz model (Martin et al. – paper is in `Readings` directory) is a phenomenological alternative:

$$B(T) = B_{max} \cdot e^{\left[ -E \cdot (T - T_{pk})^2 - e^{[E_D \cdot (T - T_{pk}) - \theta]} \right]} \quad (7.2)$$

where  $B_{max}$  is the maximum value of the available measurements and  $\theta$  a quantity that improves the fit, without any biological interpretation. The temperature parameters ( $T, T_{pk}$ ) of the Gaussian-Gompertz model (Equation 7.2) are in  ${}^\circ\text{C}$ .

### 7.3.4 Fitting data to the models

You will use Nonlinear Least Squares (NLLS) to fit the two alternative models to data, followed by model selection with AIC and BIC (also known as the Schwartz Criterion — recall the Johnson and Omland paper).

### 7.3.5 The Workflow

You will build a workflow that starts with the data and ends with a report written in `LATEX`. I suggest the following components and sequence in your workflow (you can choose to do it differently!):

1. An R script that imports the data and prepares it for NLLS fitting, with the following features:
  - It should create unique ids so that you can identify unique thermal responses (what does this mean?)

- It should filter out datasets with less than 5 data points (why?)
  - It should deal with negative and zero trait values (why?)
  - The script should add columns containing starting values of the model parameters for the NLLS fitting
  - Save the modified data to a new csv file
2. A python script that opens the new modified dataset (from step 1) and does the NLLS fitting, with the following features:
- Uses lmfit — look up submodules minimize, Parameters, Parameter, and report\_fit.
- Have a look through <http://lmfit.github.io/lmfit-py>, especially <http://lmfit.github.io/lmfit-py/fitting.html#minimize>*
- You will have to install lmfit using pip or easy\_install (you may get some errors while installing but nothing un-solvable!) — One of us can help you out if needed. *Also have a look (and try out) the examples at: <http://lmfit.github.io/lmfit-py/intro.html> & <http://lmfit.github.io/lmfit-py/parameters.html>, and whatever other examples you can find.*
- Will use the try construct because not all runs will converge. Recall the try example from R
  - Will calculate AIC, BIC,  $R^2$ , and other statistical measures of fit (you decide what you want to include)
  - Will export the results to a csv that the plotting R script (next item) can read.
3. A R script that imports the results from the previous step and plots every thermal response with both models (or none, if nothing converges) overlaid — all plots should be saved in a single separate sub-directory. Use ggplot for pretty results!
4. A L<sup>A</sup>T<sub>E</sub>X source code that generates your report.
5. A python script (saved in Code) called run\_LongPrac.py that should run the whole project, right down to compilation of the L<sup>A</sup>T<sub>E</sub>X document.

Doing all this may seem a bit scary at the start. However, you need to approach the problem systematically and methodically, and you will be OK. I suggest the following to get you started:

- Explore the data in R and get a preliminary version of the plotting script without the fitted models overlaid worked out. That will also give you a feel for the data.
- Explore the two models — be able to plot them. Write them as functions in your python script, because that's where you will use them (step 2 above) (you can use matplotlib for quick and dirty plotting and then suppress those code lines later).
- Figure out, using a minimal example (say, with one, “nice-looking” thermal response dataset) to see how the python lmfit module works. Dimitris can help you work out the minimal example, including the usage of try to catch errors in case the fitting doesn't converge.
- One thing to note is that you will need to do the NLLS fitting on the logarithm of the the function to facilitate convergence — please ask me or Dimitris if you need help on this.

## The Report

The report should,

- be written in L<sup>A</sup>T<sub>E</sub>X using the article document class

- be double-spaced, with *continuous* line numbers
- have a title, author(s) name(s) with affiliation, brief introduction with objectives of the study, and appropriate additional sections such as methods, data, results, discussion, etc — multiple authors if it is an MSc–MRes collaboration.
- must contain  $\leq 4000$  words, including references! — there should be a wordcount at the beginning of the document.
- all references should be properly cited using bibtex.
- if there are multiple authors, the contributions should be stated in detail (who did what part of the coding, etc.). I am assuming the MSc lead author will have done most of the work!

### 7.3.6 Patching together your workflow components

Use `python` for this. You can call and run all the components of the above suggested workflow, including compilation of the `LATEX`document. Look back at the notes to see how you would run these different components — we have already covered how to run `R` and compile `LATEX`using `subprocess`.

### 7.3.7 Submission

Commit and push all your work to your bitbucket repository using a directory called `CMEELongPrac` at the same level as the `Week1`, `Week2` etc. directories.

At this stage, I am not going to tell you how to organize your project — that's one of my marking criteria (see next section).

#### Marking criteria

I will be looking for the following in your submission:

- A well-organized project where code, results, data, etc., are easy to locate and inspect
- A script called `run_LongPrac.py` that runs the project successfully
- A report that contains all the components indicated above in the subsection called “The Report” — I will be looking for some original thought and synthesis in the intro and discussion
- The more thermal response curves you are able to fit, the better — that is part of the challenge!
- Quality of the presentation of the graphics and tables in your report, as well as the plots showing the fits.

#### Deadline

Commit and push your work to your git repository by xx, 5PM.

## 7.4 Errors in your `python` code

What do you want from your code? Rank the following by importance:

1. it gives me the right answer

2. it is very fast
3. it is possible to test it
4. it is easy to read
5. it uses lots of 'clever' programming techniques
6. it has lots of tests
7. it uses every language feature that you know about

Then, think about this:

- If you are very lucky, your program will crash when you run it
- If you are lucky, you will get an answer that is obviously wrong
- If you are unlucky, you won't notice until after publication
- If you are very unlucky, someone else will notice it after publication

Ultimately, most of your time could well be spent error-checking and “debugging”, not writing code.

#### 7.4.1 Unit testing

Unit testing prevents the most common mistakes and yields reliable code. Indeed, there are many reasons for testing:

- Can you prove (to yourself) that your code does what you think it does?
- Did you think about the things that might go wrong?
- Can you prove to other people that your code works?
- Does it still all still work if you fix a bug?
- Does it still all still work if you add a feature?
- Does it work with that new dataset?
- Does it work on the latest version of R?
- Does it work on Mac, Linux, Windows?
- 64 bit and 32 bit?
- Does it work on an old version of a Mac
- Does it work on Harvey, or Imperial's Linux cluster?

The idea is to write *independent* tests for the *smallest units* of code. Why the smallest units? — to be able to retain the tests upon code modification.

Let's try `doctest`, the simplest testing tool in python: simpletests for each function are embedded in the docstring. Copy the file `control_flow.py` into the file `test_control_flow.py` and edit the original function so:

```
#!/usr/bin/python

"""Some functions exemplifying the use of control statements"""

__author__ = 'Samraat Pawar (s.pawar@imperial.ac.uk)'
__version__ = '0.0.1'

import sys
import doctest # Import the doctest module
```

```

def even_or_odd(x=0):
    """Find whether a number x is even or odd.

    >>> even_or_odd(10)
    '10 is Even!'

    >>> even_or_odd(5)
    '5 is Odd!'

    whenever a float is provided, then the closest integer is used:
    >>> even_or_odd(3.2)
    '3 is Odd!'

    in case of negative numbers, the positive is taken:
    >>> even_or_odd(-2)
    '-2 is Even!'

    """
    #Define function to be tested
    if x % 2 == 0:
        return "%d is Even!" % x
    return "%d is Odd!" % x

## I SUPPRESSED THIS BLOCK: WHY?

# def main(argv):
#     # print even_or_odd(22)
#     # print even_or_odd(33)
#     # return 0

# if (__name__ == "__main__"):
#     # status = main(sys.argv)

doctest.testmod()      # To run with embedded tests

```

Now type `run test_control_flow.py -v`:

```

In []: run test_control_flow.py -v
Trying:
    even_or_odd(10)
Expecting:
    '10 is Even!'
ok
Trying:
    even_or_odd(5)
Expecting:
    '5 is Odd!'
ok
Trying:
    even_or_odd(3.2)
Expecting:
    '3 is Odd!'
ok
Trying:
    even_or_odd(-2)
Expecting:
    '-2 is Even!'
ok

```

```
1 items had no tests:
    __main__
1 items passed all tests:
    4 tests in __main__.even_or_odd
4 tests in 2 items.
4 passed and 0 failed.
Test passed.
```

You can also run doctest “on the fly”, without writing doctest.testmod() in the code by typing in a terminal: python -m doctest -v your\_function\_to\_test.py

For more complex testing, see documentation of doctest at [www.docs.python.org/2/library/doctest.html](http://www.docs.python.org/2/library/doctest.html), the package nose and the package unittest

Please start testing as early as possible, but don’t try to test everything either! Remember, it is easier to test if code is compartmentalized into functions (*why?*).

### 7.4.2 Debugging

OK, so you unit-tested, let’s go look at life through beer-goggles... BUT NO! YOU WILL RUN INTO BUGS!

Bugs happen, inevitably, in life and programming. You need to find and debug them. Banish all thoughts of littering your code with print statements to find bugs.

Enter the debugger. The command pdb turns on the python debugger. Type the following in a file and save as debugme.py in your Code directory:

```
def createabug(x):
    y = x**4
    z = 0.
    import pdb; pdb.set_trace()
    y = y/z
    return y

createabug(25)
```

Now run it:

```
In []: %run debugme.py
[lots of text]
createabug(x)
      2      y = x**4
      3      z = 0.
----> 4      y = y/z
      5      return y
      6

ZeroDivisionError: float division by zero
```

OK, so let’s %pdb it

```
In []: %pdb
```

```
Automatic pdb calling has been turned ON

In []: run debugme.py
[lots of text]
ZeroDivisionError: float division by zero
> createabug()
    3      z = 0.
----> 4      y = y/z
    5      return y

ipdb>
```

Now we're in the debugger shell, and can use the following commands to navigate and test the code line by line or block by block:

n	move to the next line
ENTER	repeat the previous command
s	“step” into function or procedure (i.e., continue the debugging inside the function, as opposed to simply run it)
p x	print variable x
c	continue until next break-point
q	quit
l	print the code surrounding the current position (you can specify how many)
r	continue until the end of the function

So let's continue our debugging:

```
ipdb> p x
25
ipdb> p y
390625
ipdb> p z
0.0
ipdb> p y/z
*** ZeroDivisionError: ZeroDivisionError
('float division by zero',)
ipdb> l
  1 def createabug(x):
  2     y = x**4
  3     z = 0.
----> 4     y = y/z
  5     return y
  6
  7 createabug(25)

ipdb> q

In []: %pdb
Automatic pdb calling has been turned OFF
```

### 7.4.3 Debugging with breakpoints

You may want to pause the program run and inspect a given line or block of code (*why?* — impromptu unit-testing is one reason). To do so, simply put this snippet of code where you want to pause and start a debugging session and then run the program again:

```
import ipdb; ipdb.set_trace()
```

Or, you can use `import pdb; pdb.set_trace()`

Alternatively, running the code with the flag `%run -d` starts a debugging session from the first line of your code (you can also specify the line to stop at). If you are serious about programming, please start using a debugger (R, Python, whatever...)!

## 7.5 Profiling in Python

Donald Knuth says: *Premature optimization is the root of all evil.* Indeed, computational speed may not be your initial concern. Also, you should focus on developing clean, reliable, reusable code rather than worrying first about how fast your code runs. However, speed will become an issue when and if your analysis or modeling becomes complex enough (e.g., food web or large network simulations). In that case, knowing which parts of your code take the most time is useful — optimizing those parts may save you lots of time

To find out what is slowing down your code you need to use “profiling”. Profiling is easy in `ipython` — simply type the magic command `%run -p your_function_name`. Let’s write a simple illustrative program and name it `profileme.py`:

```
def a_useless_function(x):
    y = 0
    # eight zeros!
    for i in range(100000000):
        y = y + i
    return 0

def a_less_useless_function(x):
    y = 0
    # five zeros!
    for i in range(100000):
        y = y + i
    return 0

def some_function(x):
    print x
    a_useless_function(x)
    a_less_useless_function(x)
    return 0

some_function(1000)
```

Now type `%run -p profileme.py`, and you should see something like:

```

54 function calls in 3.652 seconds

Ordered by: internal time

ncalls  tottime  percall  cumtime  percall filename:lineno(function)
      1    2.744    2.744    3.648    3.648 profileme.py:1(←
          a_useless_function)
      2    0.905    0.452    0.905    0.452 {range}
      1    0.002    0.002    0.003    0.003 profileme.py:8(←
          a_less_useless_function)
[more output]

```

The function `range` is taking long – we should use `xrange` instead. When iterating over a large number of values, `xrange`, unlike `range`, does not create all the values before iteration, but creates them “on demand”. E.g., `range(1000000)` yields a 4Mb+ list. So let’s modify the script:

```

def a_useless_function(x):
    y = 0
    # eight zeros!
    for i in xrange(100000000):
        y = y + i
    return 0

def a_less_useless_function(x):
    y = 0
    # five zeros!
    for i in xrange(100000):
        y = y + i
    return 0

def some_function(x):
    print x
    a_useless_function(x)
    a_less_useless_function(x)
    return 0

some_function(1000)

```

Again running the magic command `%run -p` yields:

```

52 function calls in 2.153 seconds

Ordered by: internal time

ncalls  tottime  percall  cumtime  percall filename:lineno(function)
      1    2.150    2.150    2.150    2.150 profileme2.py:1(←
          a_useless_function)
      1    0.002    0.002    0.002    0.002 profileme2.py:8(←
          a_less_useless_function)
      1    0.001    0.001    2.153    2.153 {execfile}
[more output]

```

So we saved 1.499 s! (not enough to grab a pint, but ah well...)

## 7.6 Practical

As always, add and commit all your new (functional) code and data to the version control repository: `LV1.py`, `debugme.py`, `profileme.py`.

### 7.6.1 Lotka-Volterra model problem

Copy and modify `LV1.py` into another script called `LV2.py` that does the following:

1. Take arguments for the four LV model parameters  $r$ ,  $a$ ,  $z$ ,  $e$  from the commandline

```
LV2.py arg1 arg2 ... etc
```

2. Runs the Lotka-Volterra model with prey density dependence  $rR(1 - \frac{R}{K})$ , which changes the coupled ODEs to,

$$\frac{dR}{dt} = rR\left(1 - \frac{R}{K}\right) - aCR \quad (7.3)$$

$$\frac{dC}{dt} = -zC + eaCR \quad (7.4)$$

3. Saves the plot as `.pdf` in an external results directory (`Week2/Results`)
4. The chosen parameter values should show in the plot (e.g.,  $r = 1, a = .5$ , etc)

You can change time length  $t$  too. Also include a script called `run_LV2.py` in `Code` that will run `LV2.py` with appropriate arguments.

Extra credit if you also choose appropriate values for the paramaters such that both predator and prey persist with prey density dependence.

Extra-extra credit if you can write a recursion version of the model in discrete time – it should do everything that `run_LV2.py` does. The discrete-time equation is:

$$R_{t+1} = R_t \left(1 + r \left(1 - \frac{R_t}{K}\right)\right) \quad (7.5)$$

$$C_{t+1} = C_t (1 - z + eaR_t) \quad (7.6)$$

Extra-extra-extra credit if you can write a version of the discrete-time model (eqn 7.5) simulation with a random gaussian random fluctuation in resource's growth rate at each time-step:

$$R_{t+1} = R_t \left(1 + (r + \varepsilon) \left(1 - \frac{R_t}{K}\right)\right) \quad (7.7)$$

$$C_{t+1} = C_t (1 - z + eaR_t) \quad (7.8)$$

where  $\varepsilon$  is a random fluctuation drawn from a gaussian distribution (use `scipy.stats`).

You can also add fluctuations to both populations simultaneously this way:

$$R_{t+1} = R_t \left(1 + \varepsilon + r + \left(1 - \frac{R_t}{K}\right)\right) \quad (7.9)$$

$$C_{t+1} = C_t (1 - z + \varepsilon + eaR_t) \quad (7.10)$$

If you feel up to it, you could compare the results from eqns 7.7 and 7.9.

### 7.6.2 Missing oaks problem!

1. Open and run the code `test_oaks.py` — there's a bug, for no oaks are being found! (`where's TestOaksData.csv?`)
2. Fix the bug (hint: `import ipdb; ipdb.set_trace()`)
3. Now, write doctests to make sure that, bug or no bug, your `is_an_oak` function is working as expected (hint: `>>> is_an_oak('Fagus sylvatica')` should return `False`)
4. If you wrote good doctests, you will note that you found another error that you might not have come across just by debugging (hint: what happens if you try the doctest with 'Quercuss' instead of 'Quercus'?). How would you fix the new error you found using the doctest?

## 7.7 Databases and python

Many of you will deal with complex data — and often, lots of it. Ecological and Evolutionary data are particularly complex because they contain large numbers of attributes, often measured in very different scales and units for individual taxa, populations, etc. In this scenario, storing the data in a database makes a lot of sense! You can easily include the database in your analysis workflow — indeed, that's why people use databases.

A *relational* database is a collection of interlinked (*related*) tables that altogether store a complex dataset in a logical, computer-readable format. Dividing a dataset into multiple tables minimizes redundancies. For example, if your data were sampled from three sites — then, rather than repeating the site name and description in each row in a text file, you could just specify a numerical “key” that directs to another table containing the sampling site name and description.

Finally, if you have many rows in your data file, the type of sequential access we have been using in our python and R scripts is inefficient — you should be able to instantly access any row regardless of its position

Data columns in a database are usually called *fields*, while the rows are the *records*. Here are a few things to keep in mind about databases:

- Each field typically contains only one data type (e.g., integers, floats, strings)
- Each record is a “data point”, composed of different values, one for each field — somewhat like a python tuple
- Some fields are special, and are called *keys*:
  - The *primary key* uniquely defines a record in a table (e.g., each row is identified by a unique number)
  - To allow fast retrieval, some fields (and typically all the keys) are indexed — a copy of certain columns that can be searched very efficiently
  - *Foreign keys* are keys in a table that are primary keys in another table and define relationships between the tables
- The key to designing a database is to minimize redundancy and dependency without losing the logical consistency of tables — this is called *normalization* (arguably more of an art than a science!)

Let's look at a simple example. Imagine you recorded body sizes of species from different field sites in a text file with fields:

ID	Unique ID for the record
SiteName	Name of the site
SiteLong	Longitude of the site
SiteLat	Latitude of the site
SamplingDate	Date of the sample
SamplingHour	Hour of the sampling
SamplingAvgTemp	Average air temperature on the sampling day
SamplingWaterTemp	Temperature of the water
SamplingPH	PH of the water
SpeciesCommonName	Species of the sampled individual
SpeciesLatinBinom	Latin binomial of the species
BodySize	Width of the individual
BodyWeight	Weight of the individual

It would be logical to divide the data into four tables:

*Site table:*

SiteID	ID for the site
SiteName	Name of the site
SiteLong	Longitude of the site
SiteLat	Latitude of the site

*Sample table:*

SamplingID	ID for the sampling date
SamplingDate	Date of the sample
SamplingHour	Hour of the sample
SamplingAvgTemp	Average air temperature
SamplingWaterTemp	Temperature of the water
SamplingPH	PH of the water

*Species table:*

SpeciesID	ID for the species
SpeciesCommonName	Species name
SpeciesLatinBinom	Latin binomial of the species

*Individual table:*

IndividualID	ID for the individual sampled
SpeciesID	ID for the species
SamplingID	ID for the sampling day
SiteID	ID for the site
BodySize	Width of the individual
BodyWeight	Weight of the individual

In each table, the first ID field is the primary key. The last table contains three foreign keys because each individual is associated with one species, one sampling day and one sampling site.

These structural features of a database are called its *schema*.

### 7.7.1 python with SQLite

SQLite is a simple (and very popular) SQL (Structured Query Language)-based solution for managing localized, personal databases. I can safely bet that most, if not all of you unknowingly (or knowingly!) use SQLite — it is used by MacOSX, Firefox, Acrobat Reader, iTunes, Skype, iPhone, etc.

We can easily use SQLite through Python scripting. First, install SQLite by typing in the Ubuntu terminal:

```
$ sudo apt-get install sqlite3 libsqlite3-dev
```

Also, make sure that you have the necessary package for python by typing `import sqlite3` in the python or ipython shell. Finally, you may install a GUI for SQLite3 :

```
$ sudo apt-get install sqliteman
```

Now type `sqlite3` in the Ubuntu terminal to check if SQLite successfully launches.

SQLite has very few data types (and lacks a boolean and a date type):

NULL	The value is a NULL value
INTEGER	The value is a signed integer, stored in up to or 8 bytes
REAL	The value is a floating point value, stored as in 8 bytes
TEXT	The value is a text string
BLOB	The value is a blob of data, stored exactly as it was input (useful for binary types, such as bitmap images or pdfs)

Typically, you will build a database by importing csv data — be aware that:

- Headers: the csv should have no headers
- Separators: if the comma is the separator, each record should not contain any other commas
- Quotes: there should be no quotes in the data
- Newlines: there should be no newlines

Now build your first database in SQLite! We will build it from a global dataset on wood density from Chave *et al.* (2009), Ecology Letters (should be in your Data directory). Navigate to your Code directory and in an Ubuntu terminal type:

```
$ sqlite3 ../Data/wood.db
SQLite version 3.7.9
Enter ".help" for instructions
Enter SQL statements terminated with a ";"
```

This creates an empty database in your Data directory. Now, you need to create a table with some fields:

```
sqlite> CREATE TABLE woodb (Number integer primary key,
...> Family text,
...> Binomial text,
```

```
...> WoodDensity real,
...> Region text,
...> ReferenceNumber integer);
```

Note that I am writing all SQL commands in upper case, but it is not necessary. I am using upper case here because SQL syntax is long and clunky, and it quickly becomes hard to spot (and edit) commands in long strings of complex queries.

Now let's import the dataset:

```
sqlite> .mode csv
sqlite> .import wood.csv woodb
```

So we built a table and imported a csv file into it. Now we can ask SQLite to show all the tables we currently have:

```
sqlite> .tables
woodb
```

Let's run our first *Query* (note that you need a semicolon to end a command):

```
sqlite> SELECT * FROM woodb LIMIT 5;
1,Fabaceae,"Abarema jupunba",0.78,"South America...
2,Fabaceae,"Abarema jupunba",0.66,"South America...
3,Fabaceae,"Abarema jupunba",0.551,"South America...
4,Fabaceae,"Abarema jupunba",0.534,"South America...
5,Fabaceae,"Abarema jupunba",0.551,"South America...
```

Let's turn on some nicer formatting:

```
sqlite> .mode column
sqlite> .header ON
sqlite> SELECT * FROM woodb LIMIT 5;
Number      Family      Binomial      WoodDensity ...
-----      -----      -----
1           Fabaceae   Abarema jupunba  0.78        ...
2           Fabaceae   Abarema jupunba  0.66        ...
3           Fabaceae   Abarema jupunba  0.551       ...
4           Fabaceae   Abarema jupunba  0.534       ...
5           Fabaceae   Abarema jupunba  0.551       ...
```

The main statement to select records from a table is **SELECT**:

```
sqlite> .width 40 ## NOTE: Control the width
sqlite> SELECT DISTINCT Region FROM woodb;
```

```
Region
-----
Africa (extratropical)
Africa (tropical)
Australia
Australia/PNG (tropical)
....  
  
sqlite> SELECT Binomial FROM woodb
...> WHERE Region = "Africa (tropical)";  
  
Binomial
-----
Acacia holosericea
Adansonia digitata
Afzelia africana
Afzelia africana
Afzelia africana
....  
  
sqlite> SELECT COUNT (*) FROM woodb;  
  
COUNT (*)
-----
16468  
  
sqlite> SELECT Region, COUNT(Binomial)
...> FROM woodb GROUP BY Region;  
  
Region          COUNT(Binomial)
-----  
Africa (extratropica 351
Africa (tropical) 2482
Australia        678
Australia/PNG (tropi 1560
Central America (tro 420
....  
  
sqlite> SELECT COUNT(DISTINCT Family)
...> FROM woodb;  
  
COUNT(DISTINCT Family)
-----
191  
  
sqlite> SELECT COUNT(DISTINCT Family)
...> AS FamCount
...> FROM woodb;  
  
FamCount
-----
191  
  
sqlite> SELECT Region,
...> COUNT(DISTINCT Family) AS FC
...> FROM woodb GROUP BY Region;  
  
Region          FC
-----  
Africa (extratropica 74
```

```
Africa (tropical)      66
Australia             49
Australia/PNG (tropi  91
Central America (tro  68
.....
sqlite> SELECT * # WHAT TO SELECT
...> FROM woodb # FROM WHERE
...> WHERE Region = "India" # CONDITIONS
...> AND Family = "Pinaceae";

Number          Family        Binomial       WoodDensity
-----          -----        -----
29              Pinaceae     Abies pindrow  0.38
3287             Pinaceae    Cedrus deodar   0.47
12057             Pinaceae   Picea morinda  0.4
12126             Pinaceae   Pinus longifo  0.48
12167             Pinaceae   Pinus wallich  0.43
15858             Pinaceae   Tsuga brunoni 0.38
```

The structure of the SELECT command is as follows (*Note: all characters are case insensitive*):

```
SELECT [DISTINCT] field
FROM table
WHERE predicate
GROUP BY field
HAVING predicate
ORDER BY field
LIMIT number
;
```

Let's try some more elaborate queries:

```
sqlite> SELECT Number FROM woodb LIMIT 5;

Number
-----
1
2
3
4
5

sqlite> SELECT Number
...> FROM woodb
...> WHERE Number > 100
...> AND Number < 105;

Number
-----
101
102
103
104

sqlite> SELECT Number
...> FROM woodb
...> WHERE Region = "India"
```

```
...> AND Number > 700
...> AND Number < 800;
```

Number

---

716

You can also match records using something like regular expressions. In SQL, when we use the command `LIKE`, the percent % symbol matches any sequence of zero or more characters and the underscore matches any single character. Similarly, `GLOB` uses the asterisk and the underscore.

```
sqlite> SELECT DISTINCT Region
...> FROM woodb
...> WHERE Region LIKE "_ndia";

Region
-----
India

sqlite> SELECT DISTINCT Region
...> FROM woodb
...> WHERE Region LIKE "Africa%";

Region
-----
Africa (extratropical)
Africa (tropical)

sqlite> SELECT DISTINCT Region
...> FROM woodb
...> WHERE Region GLOB "Africa*";

Region
-----
Africa (extratropica
Africa (tropical)

# NOTE THAT GLOB IS CASE SENSITIVE, WHILE LIKE IS NOT

sqlite> SELECT DISTINCT Region
...> FROM woodb
...> WHERE Region LIKE "africa%";

Region
-----
Africa (extratropical)
Africa (tropical)
```

We can also order by any column:

```
sqlite> SELECT Binomial, Family FROM
...> woodb LIMIT 5;

Binomial          Family
-----  -----
Abarema jupunba  Fabaceae
Abarema jupunba  Fabaceae
```

```

Abarema jupunba      Fabaceae
Abarema jupunba      Fabaceae
Abarema jupunba      Fabaceae

sqlite> SELECT Binomial, Family FROM
...> woodb ORDER BY Family LIMIT 5;

Binomial            Family
-----
Avicennia alba     Acanthaceae
Avicennia alba     Acanthaceae
Avicennia alba     Acanthaceae
Avicennia germinans Acanthaceae
Avicennia germinans Acanthaceae

```

So on and so forth (joining tables etc. would come next...). You can store your queries and database management commands in an .sql file (geany will take care of syntax highlighting etc.)

It is easy to access, update and manage SQLite databases with python (you should have this script file in your Code directory):

```

# import the sqlite3 library
import sqlite3

# create a connection to the database
conn = sqlite3.connect('../Data/test.db')

# to execute commands, create a "cursor"
c = conn.cursor()

# use the cursor to execute the queries
# use the triple single quote to write
# queries on several lines
c.execute('''CREATE TABLE Test
            (ID INTEGER PRIMARY KEY,
             MyVal1 INTEGER,
             MyVal2 TEXT)''')

#~c.execute('''DROP TABLE test''')

# insert the records. note that because
# we set the primary key, it will auto-increment
# therefore, set it to NULL
c.execute('''INSERT INTO Test VALUES
            (NULL, 3, 'mickey')''')

c.execute('''INSERT INTO Test VALUES
            (NULL, 4, 'mouse')''')

# when you "commit", all the commands will
# be executed
conn.commit()

# now we select the records
c.execute("SELECT * FROM TEST")

# access the next record:
print c.fetchone()

```

```

print c.fetchone()

# let's get all the records at once
c.execute("SELECT * FROM TEST")
print c.fetchall()

# insert many records at once:
# create a list of tuples
manyrecs = [(5, 'goofy'),
             (6, 'donald'),
             (7, 'duck')]

# now call executemany
c.executemany(''INSERT INTO test
                  VALUES(NULL, ?, ?)'', manyrecs)

# and commit
conn.commit()

# now let's fetch the records
# we can use the query as an iterator!
for row in c.execute('SELECT * FROM test'):
    print 'Val', row[1], 'Name', row[2]

# close the connection before exiting
conn.close()

```

You can create a database in memory, without using the disk — thus you can create and discard an SQLite database within your workflow!:

```

import sqlite3

conn = sqlite3.connect(":memory:")

c = conn.cursor()

c.execute("CREATE TABLE tt (Val TEXT)")

conn.commit()

z = [('a',), ('ab',), ('abc',), ('b',), ('c',)]

c.executemany("INSERT INTO tt VALUES (?)", z)

conn.commit()

c.execute("SELECT * FROM tt WHERE Val LIKE 'a%'").fetchall()

conn.close()

```

## 7.8 Practical

Bring today's work under version control: TestR.py, TestR.R, woodb, pythonsql.py, sqlitemem.py

### 7.8.1 Using os problem 1

Open `using_os.py` and complete the tasks assigned (hint: you might want to look at `subprocess.os.walk()`)

### 7.8.2 Using os problem 2

Open `fmr.R` and work out what it does; check that you have `NagyEtAl1999.csv`. Now write python code called `run_fmr_R.py` that:

- Runs `fmr.R` to generate the desired result
- `run_fmr_R.py` should also print to the python screen whether the run was successful, and the contents of the R console output

## 7.9 Readings and Resources

- look up `docs.python.org/2/library/index.html` – Read about the packages you think will be important to you...
- “The Definitive Guide to SQLite” is a pretty complete guide to SQLite and freely available from <http://evalenzu.mat.utfsm.cl/Docencia/2012/ SQLite.pdf>
- for databases in general, try the Stanford Introduction to Databases course: <https://www.coursera.org/course/db>
- Some of you might find the python package `biopython` particularly useful — check out [www.biopython.org/wiki/Main\\_Page](http://www.biopython.org/wiki/Main_Page), and especially, the cookbook



# Chapter 8

## Introduction to R

### 8.1 Outline of the the R module

The content and structure of the chapters in the R module are geared towards the following objectives:

- Give you an introduction to R syntax and programming conventions, assuming you have never set your eyes on R or any other programming language before — we will breeze through this as you have already been introduced to R in the Stats Week!
- Teach you principles of clean and efficient programming in R, including delightful things like vectorization and debugging.
- Teach you how to generate publication quality graphics in R — publication quality is thesis quality!
- Teach you how to develop reproducible data analyses “work flows” so you (or anybody else) can run and re-run your analyses, graphics outputs and all, in R.

You will use R a lot during the rest of your courses and probably your thesis and career — the aim is to lay down the foundations for you to become very comfortable with it!

### 8.2 What is R?

R is a freely available statistical software with strong programming capabilities. R has become incredibly popular in biology due to several factors: i) many packages are available to perform all sorts of statistical and mathematical analyses; ii) it has been developed and scrutinised by top level academic statisticians; iii) it can produce beautiful, publication-quality graphics; iv) it has a very good support for matrix-algebra.

### 8.3 Would you ever need anything other than R?

Although we can technically program in R, the programming environment is not the greatest: especially the way types are managed is problematic (e.g., often your matrix will become a vector

if it has only one column/row!), and the way errors and warnings are handled and displayed (often unintelligibly!).

Nevertheless, being able to program R means you can develop and automate your statistical analyses and the generation of figures into a reproducible work flow (there's that term again!). However, if your work also includes extensive numerical simulations, manipulation of very large matrices, bioinformatics, or complex workflows including databases, you will be much better off if you *also* know another programming language that is more versatile and computationally efficient (like python, perl or C).

In particular, python is recommended, as it is reasonably efficient and has very nice and clean syntax. With something like python, You can embed your R analysis work flow inside a more complex meta-workflow, for example, one that includes interfacing with the internet, manipulating/querying databases, and compiling L<sup>A</sup>T<sub>E</sub>Xdocuments.

But for many of you, R will do the job!

## 8.4 Installing R

Linux/Ubuntu: run the following in terminal

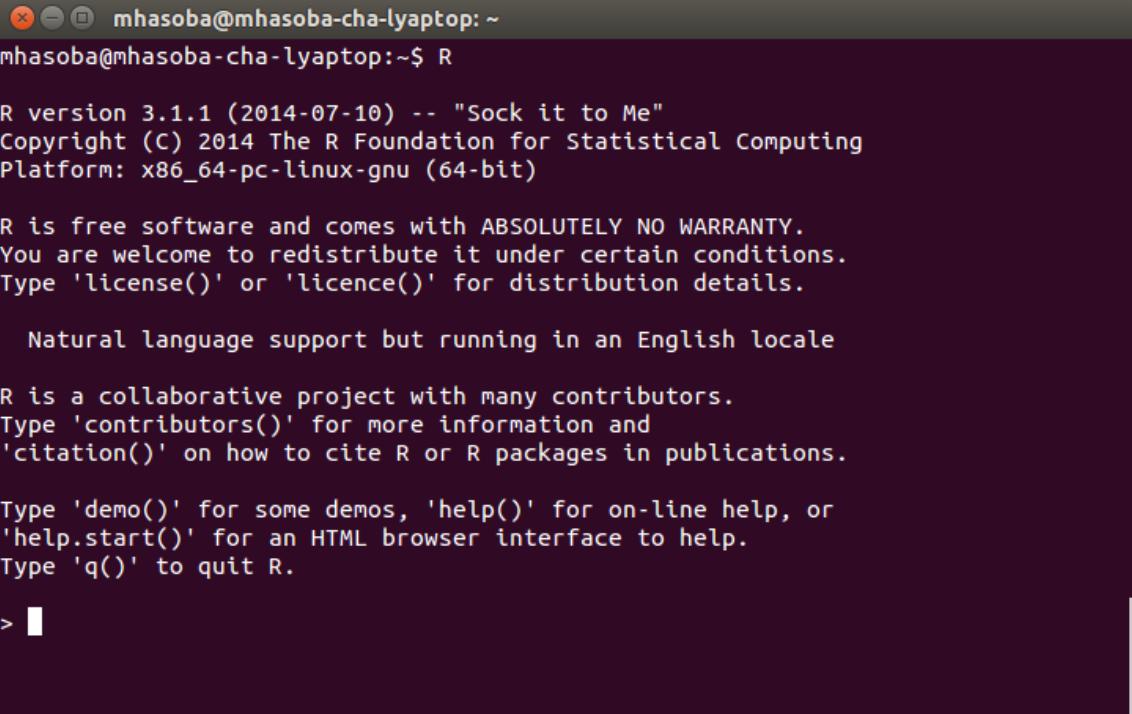
```
sudo apt-get install r-base r-base-dev
```

Mac OS X: download and install from <http://cran.r-project.org/bin/macosx/>

Windows: download and install from <http://cran.r-project.org/bin/windows/base/>

## 8.5 Getting started

Launch R (From Applications menu on Window or Mac, from terminal in Linux/Ubuntu) — it should look something like this (on Linux/Ubuntu or Mac terminal):



```
mhasoba@mhasoba-cha-lyaptop:~$ R
R version 3.1.1 (2014-07-10) -- "Sock it to Me"
Copyright (C) 2014 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

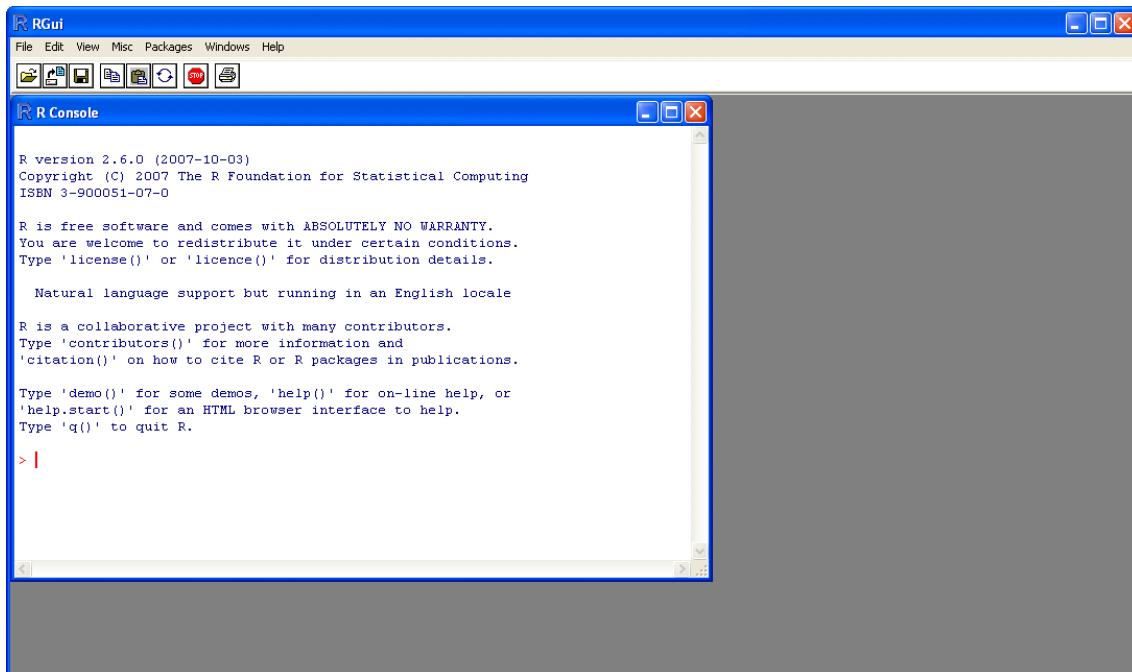
Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> |
```

Or like this (Windows “console”, similar in Mac):



You can also use an IDE (Interactive Development Environment) that can offer delights like syntax highlighting (google it!), such as RStudio, geany, vim, etc.

## 8.6 Useful R commands

<code>ls()</code>	list all the variables in the work space
<code>rm('a', 'b')</code>	remove variable(s) a and b
<code>rm(list=ls())</code>	remove all variable(s)
<code>getwd()</code>	get current working directory
<code>setwd('Path')</code>	set working directory to Path
<code>q()</code>	quit R
<code>?Command</code>	show the documentation of Command
<code>??Keyword</code>	search the all packages/functions with Keyword, “fuzzy search”

## 8.7 Some Basics

Try out the following in the R console:

```
> a <- 4 # assignment
> a
[1] 4
> a*a # product
[1] 16
> a_squared <- a*a
> sqrt(a_squared) # square root
[1] 4
> v <- c(0, 1, 2, 3, 4) # c: "concatenate"
```

`c()` (concatenate) is one of the most commonly used functions — Don’t forget it! (try `?c`)

Note that any text after a “#” is ignored by R — handy for commenting. In general, please comment your code and scripts, for *everybody*’s sake!

```
> v # Display the vector variable you created
[1] 0 1 2 3 4
> is.vector(v) # check if it's a vector
[1] TRUE
> mean(v) # mean
[1] 2
> var(v) # variance
[1] 2.5
> median(v) # median
[1] 2
> sum(v) # sum all elements
[1] 10
> prod(v + 1) # multiply
[1] 120
> length(v) # length of vector
[1] 5
```

### 8.7.1 Variable names and Tabbing

```
> wing.width.cm <- 1.2 #Using dot notation
> wing.length.cm <- c(4.7, 5.2, 4.8)
```

Using tabbing with dot notation can be handy. Type:

```
> wing.
```

And then hit the `tab` key. This is handy, but good style and readability is more important than just convenient variable names. Variable names should be as obvious as possible, not over-long!

### 8.7.2 R likes E Notation

```
> 1E4
[1] 10000
> 1e4
[1] 10000
> 5e-2
[1] 0.05
```

R uses *E* notation to print very large or small numbers:

```
> 1E4 ^ 2
[1] 1e+08
> 1 / 3 / 1e8
[1] 3.333333e-09
```

### 8.7.3 Operators

The usual operators are available in R (slight differences from python):

<code>+</code>	Addition
<code>-</code>	Subtraction
<code>*</code>	Multiplication
<code>/</code>	Division
<code>^</code>	Power
<code>%%</code>	Modulo
<code>%/%</code>	Integer division
<code>==</code>	Equals
<code>!=</code>	Differs
<code>&gt;</code>	Greater
<code>&gt;=</code>	Greater or equal
<code>&amp;</code>	Logical and
<code> </code>	Logical or
<code>!</code>	Logical not

### 8.7.4 When things go wrong

Syntax errors are those where you've just made a typing mistake. Here are some common problems:

- missing close bracket leads to continuation line.

```
> x <- (1 + (2 * 3)
+
```

Hit Ctrl C (see below) or keep typing!

- Too many parentheses: `2 + (2 * 3)`
- wrong/mismatched brackets (see next subsection).
- Do not mix double quotes and single quotes.
- When things seem to take too long, try `Ctrl + C`

### 8.7.5 Types of parentheses

R has a somewhat confusing array of parentheses that you need to get used to:

- `f(3, 4)` – call the function `f`, with `arg1=3, arg2=4`.
- `a + (b*c)` – use to enforce order over which statements are executed.
- `{ expr1; expr2; ...exprn }` – group a set of expressions into one compound expression. Value returned is value of last expression; used in looping/conditionals.
- `x[4]` – get the 4th element of the vector `x`.
- `l[[3]]` – get the 3rd element of some list `l`, and return it (compare with `l[3]` which returns a list with just the 3rd element inside – more on lists in next section).

## 8.8 Data types

Like python (why python? Ask the CMEEs!), R comes with data-types. Mastering these will help you write better, more efficient programs and also handle diverse data types between datasets. Now get back into R (if you quit R using `q()`), and type:

### 8.8.1 Vectors

Vectors are a fundamental object for R. Scalars (single numbers) are treated as vector of length 1.

```
> a <- 5
> is.vector(a)
[1] TRUE
> v1 <- c(0.02, 0.5, 1)
> v2 <- c("a", "bc", "def", "ghij")
> v3 <- c(TRUE, TRUE, FALSE)
```

### 8.8.2 Matrices and arrays

R has many functions to manipulate matrices (two-dimensional vectors) and arrays (multi-dimensional vectors).

```
> m1 <- matrix(1:25, 5, 5)
> m1
[,1] [,2] [,3] [,4] [,5]
[1,]    1     6    11    16    21
[2,]    2     7    12    17    22
[3,]    3     8    13    18    23
[4,]    4     9    14    19    24
[5,]    5    10    15    20    25
> m1 <- matrix(1:25, 5, 5, byrow=TRUE)
> m1
[,1] [,2] [,3] [,4] [,5]
[1,]    1     2     3     4     5
[2,]    6     7     8     9    10
[3,]   11    12    13    14    15
[4,]   16    17    18    19    20
[5,]   21    22    23    24    25
> dim(m1)
[1] 5 5
> m1[1,2]
[1] 2
> m1[1,2:4]
[1] 2 3 4
> m1[1:2,2:4]
[,1] [,2] [,3]
[1,]    2     3     4
[2,]    7     8     9
> arr1 <- array(1:50, c(5, 5, 2))
> arr1
, , 1

[,1] [,2] [,3] [,4] [,5]
[1,]    1     6    11    16    21
[2,]    2     7    12    17    22
[3,]    3     8    13    18    23
[4,]    4     9    14    19    24
[5,]    5    10    15    20    25

, , 2

[,1] [,2] [,3] [,4] [,5]
[1,]   26    31    36    41    46
[2,]   27    32    37    42    47
[3,]   28    33    38    43    48
[4,]   29    34    39    44    49
[5,]   30    35    40    45    50
```

### 8.8.3 Data frames

This is a very important data type that is peculiar to R. It is great for storing your data. Basically, it's a two-dimensional table in which each column can contain a different data type (e.g., numbers, strings, boolean). You can think of a dataframe as a spreadsheet. Data frames are great for plotting

your data, performing regressions and such. Later (Chapter 2) you will see that fancy plotting using `ggplot` demands the use of dataframes. Now try the following:

```
> Col1 <- 1:10
> Col1
[1] 1 2 3 4 5 6 7 8 9 10
> Col2 <- LETTERS[1:10]
> Col2
[1] "A" "B" "C" "D" "E" "F" "G" "H" "I" "J"
> Col3 <- runif(10) # 10 random num. from uniform
> Col3
[1] 0.29109 0.91495 0.64962 0.95503 0.26589 0.02482 0.59718
[8] 0.99134 0.98786 0.86168
> MyDF <- data.frame(Col1, Col2, Col3)
> MyDF
   Col1 Col2     Col3
1     1    A 0.2910981
2     2    B 0.9149558
3     3    C 0.6496248
4     4    D 0.9550331
5     5    E 0.2658936
6     6    F 0.0248217
7     7    G 0.5971868
8     8    H 0.9913407
9     9    I 0.9878679
10    10   J 0.8616854
> names(MyDF) <- c("A.name", "another", "another.one")
> MyDF
  A.name another another.one
1     1        A 0.2910981
2     2        B 0.9149558
3     3        C 0.6496248
4     4        D 0.9550331
5     5        E 0.2658936
6     6        F 0.0248217
7     7        G 0.5971868
8     8        H 0.9913407
9     9        I 0.9878679
10    10       J 0.8616854
> MyDF$A.name
[1] 1 2 3 4 5 6 7 8 9 10
> MyDF[,1]
[1] 1 2 3 4 5 6 7 8 9 10
> MyDF[c("A.name", "another")]
  A.name another
1     1        A
2     2        B
3     3        C
4     4        D
5     5        E
6     6        F
7     7        G
8     8        H
9     9        I
10    10       J
> class(MyDF)
[1] "data.frame"
> str(MyDF) # a very useful command!
'data.frame': 10 obs. of 3 variables:
 $ A.name : int 1 2 3 4 5 6 7 8 9 10
```

```
$ another      : Factor w/ 10 levels "A","B","C","D",...
$ another.one: num  0.291 0.915 0.65 0.955 0.266 ...
```

### 8.8.4 Lists

A list is used to collect a group of objects of different sizes and types. A list is simply an ordered collection of objects (that can be other variables – a bit like python lists!).

```
> l1 <- list(names=c("Fred", "Bob"), ages=c(42, 77, 13, 91))
> l1
$names
[1] "Fred" "Bob"

$ages
[1] 42 77 13 91

> l1[[1]]
[1] "Fred" "Bob"
> l1[[2]]
[1] 42 77 13 91
> l1[["ages"]]
[1] 42 77 13 91
> l1$ages
[1] 42 77 13 91
```

You can build lists of lists too. Lists are often returned as the result of a complex function (e.g. linear model fitting using `lm()`) to return all relevant information in one object.

## 8.9 Variable Types, Type Conversion and Special Values

There are different kinds of data variable types such as integer, float (including real numbers), and string (e.g., text words). Beware of the difference between NA (Not Available) and NaN (Not a Number).

```
> as.integer(3.1)
[1] 3
> as.real(4)
[1] 4
> as.roman(155)
[1] CLV
> as.character(155)
[1] "155"
> as.logical(5)
[1] TRUE
> as.logical(0)
[1] FALSE
> b <- NA
> is.na(b)
[1] TRUE
> b <- 0./0.
> b
[1] NaN
```

```
> is.nan(b)
[1] TRUE
> b <- 5/0
> b
[1] Inf
> is.nan(b)
[1] FALSE
> is.infinite(b)
[1] TRUE
> is.finite(b)
[1] FALSE
> is.finite(0/0)
[1] FALSE
```

## 8.10 Creating and Manipulating Data structures

### 8.10.1 Sequences

The `:` operator creates vectors of sequential integers:

```
> years <- 1990:2009
> years
[1] 1990 1991 1992 1993 1994 1995 1996 1997 1998 1999
[11] 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009

> years <- 2009:1990 # or in reverse order
> years
[1] 2009 2008 2007 2006 2005 2004 2003 2002 2001 2000
[11] 1999 1998 1997 1996 1995 1994 1993 1992 1991 1990
```

For sequences of fractional numbers, you have to use `seq()`:

```
> seq(1, 10, 0.5)
[1] 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0 5.5 6.0 6.5 7.0 7.5 ←
     8.0
[16] 8.5 9.0 9.5 10.0
```

You can also `seq(from=1, to=10, by=0.5)` OR `seq(from=1, by=0.5, to=10)` with the same effect (try it) — this explicit, “argument matching” approach is partly why R is so popular.

### 8.10.2 Strings and Pasting

R’s string handling ain’t elegant or pretty (unlike python!), but it works:

```
> species.name <- "Quercus robur" #double quotes
> species.name
[1] "Quercus robur"
> species.name <- 'Fraxinus excelsior' #single quotes
> species.name
```

```
[1] "Fraxinus excelsior"
> paste("Quercus", "robur")
[1] "Quercus robur"
> paste("Quercus", "robur", sep = "") #Get rid of space
"Quercusrobur"
> paste("Quercus", "robur", sep = ", ") #insert comma to separate
```

And as is the case with so many R functions, pasting works on vectors:

```
> paste('Year is:', 1990:2000)
[1] "Year is: 1990" "Year is: 1991" "Year is: 1992" "Year is: 1993"
[5] "Year is: 1994" "Year is: 1995" "Year is: 1996" "Year is: 1997"
[9] "Year is: 1998" "Year is: 1999" "Year is: 2000"
```

Note that this last example creates a vector of 11 strings.

### 8.10.3 Indices and Indexing

Every element of a vector in R has an order: the first value, second, third, etc. To illustrate this, type:

```
> MyVar <- c( 'a' , 'b' , 'c' , 'd' , 'e' ) # create a simple vector
```

Then, square brackets extract values based on their position in the order:

```
> MyVar[1] # Show element in first position
[1] "a"
> MyVar[4]
[1] "d" # Show element in fourth position
```

The values in square brackets are called “indices” — they give the index (position) of the required value. We can also select sets of values in different orders, or repeat values:

```
> MyVar[c(3,2,1)] # reverse order
[1] "c" "b" "a"
MyVar[c(1,1,5,5)] # repeat indices
[1] "a" "a" "e" "e"
```

So you can manipulate vectors by indexing:

```
> v <- c(0, 1, 2, 3, 4) # Re-create the vector variable v
> v[3] # access one element
[1] 2
> v[1:3] # access sequential elements
[1] 0 1 2
> v[-3] # remove elements
[1] 0 1 3 4
> v[c(1, 4)] # access non-sequential
[1] 0 3
```

### 8.10.4 Recycling

When vectors are of different lengths, R will recycle the shorter one to make a vector of the same length:

```
a <- c(1,5) + 2
x <- c(1,2); y <- c(5,3,9,2)
x + y
x + c(y,1) ## somewhat strange!
```

Recycling is convenient, but dangerous!

### 8.10.5 Basic vector-matrix operations

```
> v2 <- v
> v2 <- v2*v2 # whole-vector operation
> v2
[1] 0 2 4 6 8
> v * v2 # product element-wise
[1] 0 2 8 18 32
> t(v) # transpose the vector
[,1] [,2] [,3] [,4] [,5]
[1,] 0 1 2 3 4
> v %*% t(v) # matrix/vector product
[,1] [,2] [,3] [,4] [,5]
[1,] 0 0 0 0 0
[2,] 0 1 2 3 4
[3,] 0 2 4 6 8
[4,] 0 3 6 9 12
[5,] 0 4 8 12 16
> v3 <- 1:7 # assign using sequence
> v3
[1] 1 2 3 4 5 6 7
> v4 <- c(v2, v3) # concatenate vectors
> v4
[1] 0 2 4 6 8 1 2 3 4 5 6 7
> q() # quit
```

## 8.11 Your R Analysis Workflow

### 8.11.1 The Working Directory

Before we go any further, let's get ourselves organized. In R, type:

```
> getwd()
```

This tells you what the current “working directory” is.

In using R for an analysis, you will likely use and create several files. This means that it is sensible to create a folder (directory) to keep all code files together. You can then set R to work from this

directory, so that files are easy to find and run — this will be your working directory. So now do the following:

- ★ Create a directory called Week5 in an appropriate location (For CMEEs, in CMEECourseWork for others, some place you can remember!)
- ★ Create subdirectories within CMEECourseWork/Week5 called Code, Data, and Results

You can create directories using `dir.create()` within R

**Use relative paths** Using relative paths in your R scripts and code will make your code computer independent and your life better! For example, in R, `../Data/mydata.txt` specifies a file named `mydata.txt` located in the “parent” of the current directory. Now, set the working directory to be `Week5/Code`:

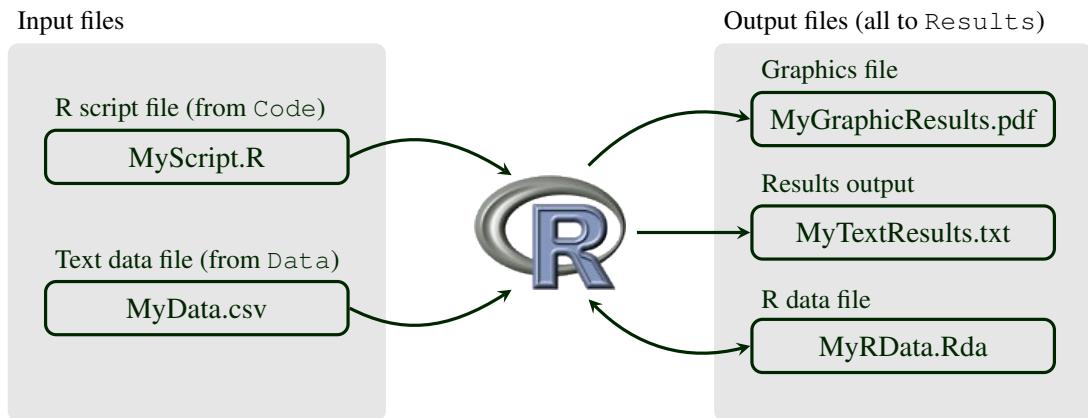
```
> setwd("FullPathUptoHere/Week5/Code")
> dir()
```

Note that `FullPathUptoHere` is not to be taken literally and entered! For example, if you created `Week5` in `H:`, you would use `setwd("H:/MyICStatsModules/Code")`

For CMEE, it would be `setwd("FullPathUptoHere/CMEECourseWork/Week5/Code")`

## 8.12 The R analysis workflow

Your typical R analysis workflow will be as follows:



Some details on each kind of file:

**R script files** These are plain text files containing all the R code needed for an analysis. These should always be created with a simple text editor like Notepad (Windows),TextEdit (MacOS) or Geany (Linux) and saved with the extension `*.R`. We will use the built-in editor in R in this class. Alternatively, RStudio is a good option because it also works across platforms. Try it. A big advantage of something like RStudio is that you will get syntax highlighting, which is very handy and will make R programming far more convenient and error-free. You should *never* use Word to save or edit these files as R can only read code from plain text files.

**Text data files** These are files of data in plain text format containing one or more columns of data (numbers, strings, or both). Although there are several format options, we will typically be

using csv files, where the entries are separated by commas. These are easy to create and export from Excel (if that's what you use...).<sup>1</sup>

**Results output files** These are a plain text files containing your results, such as the summary output of a regression or ANOVA analysis. Typically, you will output your results in a table format where the columns are separated by commas (csv) or tabs (tab-delimited).

**Graphics files** R can export graphics in a wide range of formats. This can be done automatically from R code and we will look at this later but you can also select a graphics window and click 'File > Save as...'

**Rdata files** You can save any data loaded or created in R, including model outputs and other things, into a single Rdata file. These are not plain text and can only be read by R, but can hold all the data from an analysis in a single handy location. I never use these, but you can, if you want.

## 8.13 Importing and Exporting Data

Now we are ready to see how to import and export data in R, typically the first step of your analysis. The best option is to have your data in a comma separated value text file or in a tab separated file. Then, you can use the function `read.csv` (or `read.table`) to import your data (note the relative paths!):

```
> MyData <- read.csv("../Data/trees.csv")
> head(MyData) # Have a quick look at the data frame
> str(MyData) # Have a quick look at the column types
> MyData <- read.csv("../Data/trees.csv", header = TRUE) # with headers
> MyData <- read.table("../Data/trees.csv", sep = ',', 
header = TRUE) # A more general way
> head(MyData)
> MyData <- read.csv("../Data/trees.csv", skip = 5) # skip first 5 lines
```

Note that the resulting `MyData` in your workspace is an R dataframe. You can also save your data frames using `write.table` or `write.csv`:

```
> write.csv(MyData, "../Results/MyData.csv")
> dir("../Results/") # Check if it worked
> write.table(MyData[1], file = "../Results/MyData.csv", append=TRUE) # ←
  append
> write.csv(MyData, "../Results/MyData.csv", row.names=TRUE) # write row ←
  names
> write.table(MyData, "../Results/MyData.csv", col.names=FALSE) # ignore col ←
  names
```

---

<sup>1</sup>If you are using a computer from elsewhere in the EU, Excel may use a comma ( $\pi = 3,1416$ ) instead of a decimal point ( $\pi = 3.1416$ ). In this case, csv files may use a semi-colon to separate columns and you can use the alternative function `read.csv2()` to read them into R.

## 8.14 Writing Functions

R lets you write your own functions. The syntax is quite simple, with each function accepting arguments and returning a value:

```
MyFunction <- function(Arg1, Arg2) {
  ## statements involving Arg1, Arg2
  return (ReturnValue)
}
```

- \* Type the following in a script file called `TreeHeight.R`, save it in your `Code` directory and run it using `source`:

```
# This function calculates heights of trees
# from the angle of elevation and the distance
# from the base using the trigonometric formula
# height = distance * tan(radians)
#
# Arguments:
# degrees      The angle of elevation
# distance     The distance from base
#
# Output:
# The height of the tree, same units as "distance"

TreeHeight <- function(degrees, distance)
{
  radians <- degrees * pi / 180
  height <- distance * tan(radians)
  print(paste("Tree height is:", height))
  return (height)
}

TreeHeight(37, 40)
```

## 8.15 Control statements

In R, you can write `if`, `then`, `else` statements, and `for` and `while` loops like any programming language. However, loops are slow in R, so use them sparingly.

- \* Type the following in a script file called `control.R` (save it in your `Code` directory)

```
## If statement
a <- TRUE
if (a == TRUE) {
  print ("a is TRUE")
} else {
  print ("a is FALSE")
```

```

## On a single line
z <- runif(1) ##random number
if (z <= 0.5) {
print ("Less than a quarter")}

## For loop using a sequence
for (i in 1:100){
  j <- i * i
  print(paste(i, " squared is", j))
}

## For loop over vector of strings
for(species in c('Heliodoxa rubinoides',
                 'Boissonneaua jardini',
                 'Sula nebulosus'))
{
  print(paste('The species is', species))
}

## for loop using a vector
v1 <- c("a","bc","def")
for (i in v1){
  print(i)
}

## While loop
i <- 0
while (i<100){
  i <- i+1
  print(i^2)
}

```

## 8.16 Running R code

You can run the code you wrote in blocks to test and understand it:

- ★ Place the cursor on the first line of code and run it by pressing the keyboard shortcut (PC: **ctrl+R**, Mac: **command+enter**, Linux: **ctrl+enter** if you are using **geany**).

Typing in commands interactively or running in blocks is good for starters, but you will want to switch to putting your sequence of commands into a script file, and then ask R to run those commands. This is necessary also because you will want to just run your full analysis and output all the results. The way to run **\*.R** script/code from the command line is to **source** it. This causes R to accept code input from a named file and run it:

```
> source("control.R") # Assuming you are in Code directory!
```

Note that you will need to add the directory path to the file name (**control.R** in the above example), if the file is not in your working directory. For example, **../Code/control.R** if you are in, say, **Data**.

## 8.17 Useful R Functions

There are a number of very useful functions available by default (in the “base packages”).

### 8.17.1 Mathematical

<code>log(x)</code>	Natural logarithm
<code>log10(x)</code>	Logarithm in base 10
<code>exp(x)</code>	$e^x$
<code>abs(x)</code>	Absolute value
<code>floor(x)</code>	Largest integer $< x$
<code>ceiling(x)</code>	Smallest integer $> x$
<code>pi</code>	$\pi$
<code>sqrt(x)</code>	$\sqrt{x}$
<code>sin(x)</code>	Sine function

### 8.17.2 Strings

<code>strsplit(x, ';')</code>	Split the string according to ‘;’
<code>nchar(x)</code>	Number of characters
<code>toupper(x)</code>	Set to upper case
<code>tolower(x)</code>	Set to lower case
<code>paste(x1, x2, sep=';')</code>	Join the strings inserting ‘;’

### 8.17.3 Statistical

<code>mean(x)</code>	Compute mean (of a vector or matrix)
<code>sd(x)</code>	Standard deviation
<code>var(x)</code>	Variance
<code>median(x)</code>	Median
<code>quantile(x, 0.05)</code>	Compute the 0.05 quantile
<code>range(x)</code>	Range of the data
<code>min(x)</code>	Minimum
<code>max(x)</code>	Maximum
<code>sum(x)</code>	Sum all elements

### 8.17.4 Random number distributions

<code>rnorm(10, m=0, sd=1)</code>	Draw 10 normal random numbers with mean 0 and s.d. 1
<code>dnorm(x, m=0, sd=1)</code>	Density function
<code>qnorm(x, m=0, sd=1)</code>	Cumulative density function
<code>runif(20, min=0, max=2)</code>	Twenty random numbers from uniform [0,2]
<code>rpois(20, lambda=10)</code>	Twenty random numbers from Poisson( $\lambda$ )

## 8.18 Packages

The main strength of R is that users can easily build packages and share them through `cran.r-project.org`. There are packages to do most statistical and mathematical analysis you might conceive, so check them out before reinventing the wheel! Visit `cran.r-project.org` and go to Packages to see a list and a brief description. To install a package, within R type `install.packages()` and choose the package to install.

## 8.19 Practical

Modify the script `TreeHeight.R` so that it does the following:

- ★ Loads `trees.csv` and calculates tree heights for all trees in the data. Note that the distances have been measured in meters. (Hint: use relative paths).
- ★ Creates a csv output file called `TreeHts.csv` in `Results` that contains the calculated tree heights along with the original data in the following format (only first two rows and headers shown):

```
"Species", "Distance.m", "Angle.degrees", "Tree.Height.m"
"Populus tremula", 31.6658337740228, 41.2826361937914, 25.462680727681
"Quercus robur", 45.984992608428, 44.5359166583512, 46.094124200205
```

## 8.20 Readings

- Use the internet! Google “R tutorial”, and plenty will pop up. Choose one that seems the most intuitive to you.
- The Use R! series (the yellow books) by Springer are really good. In particular, consider: “A Beginner’s Guide to R”, “R by Example”, “Numerical Ecology With R”, “ggplot2” (we’ll see this in Chapter 2), “A Primer of Ecology with R”, “Nonlinear Regression with R”, “Analysis of Phylogenetics and Evolution with R”.
- For more focus on dynamical models: Soetaert & Herman. 2009 ‘A practical guide to ecological modelling: using R as a simulation platform’.
- There are excellent websites besides cran. In particular, check out [www.statmethods.net](http://www.statmethods.net) and [http://en.wikibooks.org/wiki/R\\_Programming](http://en.wikibooks.org/wiki/R_Programming).
- For those who are coming with Matlab experience: <http://www.math.umaine.edu/~hiebeler/comp/matlabR.html>
- Bolker, B. M.: Ecological Models and Data in R (eBook and Hardcover available).
- Beckerman, A. P. & Petchey, O. L. (2012) Getting started with R: an introduction for biologists. Oxford, Oxford University Press.  
Very basic, good if you are really stuck at the outset.

- Crawley, R. (2013) *The R book*. 2nd edition. Chichester, Wiley.  
Excellent but enormous reference book, code and data available from [www.bio.ic.ac.uk/research/mjccraw/therbook/index.htm](http://www.bio.ic.ac.uk/research/mjccraw/therbook/index.htm)



# Chapter 9

## Plotting and graphics in R

### 9.1 Basic plotting

R can produce beautiful graphics, without the time consuming and fiddly methods that you might have used in Excel or equivalent (it's not you, it's Excel!). Later, you can explore how the package `ggplot2`, can allow the rapid creation of truly elegant publication-grade graphics. However, in many cases you just want to quickly plot the data for exploratory analysis of your data. Here are the basic plotting commands:

<code>plot(x, y)</code>	Scatterplot
<code>plot(y~x)</code>	Scatterplot with y as a response variable
<code>hist(mydata)</code>	Histogram
<code>barplot(mydata)</code>	Bar plot
<code>points(y1~x1)</code>	Add another series of points
<code>boxplot(y~x)</code>	Boxplot

Let's try some basic plotting. As a case study, we will use a simplified version of a dataset on predator-prey body mass ratios taken from the Ecological Archives of the ESA (Barnes *et al.* 2008, Ecology 89:881).



These data should be in your Data directory. Let's import the data:

```
> MyDF <- read.csv("../Data/EcolArchives-E089-51-D1.csv")
> dim(MyDF)
[1] 34931      15
> MyDF$Record.number           MyDF$Predator.mass
MyDF$In.refID                  MyDF$Prey
MyDF$IndividualID               MyDF$Prey.common.name
MyDF$Predator                   MyDF$Prey.taxon
MyDF$Predator.common.name       MyDF$Prey.mass
MyDF$Predator.taxon             MyDF$Prey.mass.unit
MyDF$Predator.lifestage          MyDF$Location
MyDF>Type.of.feeding.interaction
```

### 9.1.1 Scatter Plot

Let's start by plotting Predator mass vs. Prey mass:

```
> plot(MyDF$Predator.mass, MyDF$Prey.mass)
```

That doesn't look very nice! Let's try taking logarithms (why?).

```
> plot(log(MyDF$Predator.mass), log(MyDF$Prey.mass))
```

We can change almost any aspect of the resulting graph; let's change the symbols by specifying the plot characters using `pch`:

```
> plot(log(MyDF$Predator.mass),log(MyDF$Prey.mass),pch=20) # Change marker
> plot(log(MyDF$Predator.mass),log(MyDF$Prey.mass),pch=20,
      xlab = "Predator Mass (kg)", ylab = "Prey Mass (kg)") # Add labels
```

### 9.1.2 Histograms

Now plot a histogram of Predator body masses:

```
> hist(MyDF$Predator.mass)
> hist(log(MyDF$Predator.mass),
      xlab = "Predator Mass (kg)", ylab = "Count") # labels
> hist(log(MyDF$Predator.mass),xlab="Predator Mass (kg)",ylab="Count",
      col = "lightblue", border = "pink") # Change bar and borders colors
```

### 9.1.3 Subplots

We can also plot both predator and prey body masses in different sub-plots using `par`

```
> par(mfcol=c(2,1), lwd = 1.5) #initialize multi-paneled plot
> par(mfg = c(1,1))
> hist(log(MyDF$Predator.mass),
      xlab = "Predator Mass (kg)", ylab = "Count",
      col = "lightblue", border = "pink",
      main = 'Predator') # Add title
> par(mfg = c(2,1));
> hist(log(MyDF$Prey.mass),
      xlab="Prey Mass (kg)",ylab="Count",
      col = "lightgreen", border = "pink",
      main = 'prey')
```

### 9.1.4 Overlaying plots

Better still, we would like to see if the predator mass and prey mass distributions are similar by overlaying them.

```
> hist(log(MyDF$Predator.mass), # Predator histogram
      xlab="Body Mass (kg)", ylab="Count",
      col = rgb(1, 0, 0, 0.5), # Note 'rgb'
      main = "Overlay")
> hist(log(MyDF$Prey.mass), col = rgb(0, 0, 1, 0.5), add = T) # Plot prey
> legend('topleft',c('Predators','Prey'), # Add legend
        fill=c(rgb(1, 0, 0, 0.5), rgb(0, 0, 1, 0.5))) # Define legend colors
```

### 9.1.5 Exercise

We can do a lot more beautification! As an exercise, try adjusting the bin widths to make them same for the predator and prey, and making the x and y labels larger and in boldface.

### 9.1.6 Saving your graphics

And you can also save the figure in a vector graphics format (important to learn to do this!). PDF is a good option:

```
> pdf("../Results/Pred_Prey_Overlay.pdf", # Open blank pdf page
      11.7, 8.3) # These numbers are page dimensions in inches
> hist(log(MyDF$Predator.mass), # Plot predator histogram (note 'rgb')
       xlab="Body Mass (kg)", ylab="Count",
       col = rgb(1, 0, 0, 0.5),
       main = "Overlay")
> hist(log(MyDF$Prey.mass), # Plot prey weights
       col = rgb(0, 0, 1, 0.5),
       add = T) # Add to same plot = TRUE
> legend('topleft',c('Predators','Prey'), # Add legend
         fill=c(rgb(1, 0, 0, 0.5), rgb(0, 0, 1, 0.5)))
> dev.off()
```

You can also try other graphic output formats. For example, `png()` (a raster format) instead of `pdf()`.

### 9.1.7 Boxplots

Now, let's try plotting boxplots instead of histograms. These are useful for getting a visual summary of your data :

```
> boxplot(log(MyDF$Predator.mass) ~ MyDF$Location, # Why the tilde?
           xlab = "Location", ylab = "Predator Mass",
           main = "Predator mass by location")
```

That's a lot of locations! You will need an appropriately wide plot to see all the boxplots adequately. Let's try boxplots by feeding interaction type:

```
> boxplot(log(MyDF$Predator.mass) ~ MyDF>Type.of.feeding.interaction,
           xlab = "Location", ylab = "Predator Mass",
           main = "Predator mass by feeding interaction type")
```

### 9.1.8 Practical 2.1

You can also make lattice graphs to avoid the somewhat laborious `par()` approach above. For this, you will need:

```
> library(lattice)
```

A lattice plot of the above data for predator mass could look like Fig. 9.1.8 (as a density plot). This was generated using (and printing to a pdf with particular dimensions):

```
> densityplot(~log(Predator.mass) | Type.of.feeding.interaction, data=MyDF)
```

Look up <http://www.statmethods.net/advgraphs/trellis.html> and the lattice package help.

In this practical, you will write script that draws and saves three lattice graphs by feeding interaction type: one of predator mass , one of prey mass and one of the size ratio of prey mass over predator mass. Note that you would want to use logarithms of masses (or mass-ratios) for all three plots. In addition, the script will calculate the mean and median predator mass, prey mass and predator-prey size-ratios to a csv file. The workflow would be:

- \* Write a script file called PP\_Lattice.R and save it in the Code directory — sourcing or running this script should result in three files called Pred\_Lattice.pdf, Prey\_Lattice.pdf, and SizeRatio\_Lattice.pdf being saved in the Results directory (the names are self-explanatory, I hope).
- \* In addition, the script should calculate the mean log predator mass, prey mass, and predator-prey size ratio, and save it as a single csv output table called PP\_Results.csv to the Results directory. The table should have appropriate headers (e.g., Feeding type, mean, median). (Hint: you will have to initialize a new dataframe in the script to first store the calculations)
- \* The script should be self-sufficient and not need any external inputs — it should import the above predator-prey dataset from the appropriate directory, and save the graphic plots to the appropriate directory (Hint: use relative paths).

## 9.2 Publication-quality figures in R

R can produce beautiful graphics, but it takes a lot of work to obtain the desired result. This is because the starting point is pretty much a “bare” plot, and adding features commonly required for publication-grade figures (legends, statistics, regressions, etc.) can be quite involved.

Moreover, it is very difficult to switch from one representation of the data to another (i.e., from boxplots to scatterplots), or to plot several dataset together. To overcome these issues, the R package ggplot2 is very powerful. It can be used to produce truly high-quality graphics for papers, theses and reports. *One thing to note though is that at present, ggplot2 cannot be used to create 3D graphs or mosaic plots.*

ggplot2 differs from other approaches as it attempts to provide a “grammar” for graphics in which each layer is the equivalent of a verb, subject etc. and a plot is the equivalent of a sentence. All graphs start with a layer showing the data, other layers and commands are added to modify the plot.

For a complete reference, please see the book “ggplot2: Elegant Graphics for Data Analysis”, by H. Wickham. Also, the website [ggplot2.org](http://ggplot2.org) a great resource. To install ggplot2, open a session of R and type (launch R using sudo R in Linux first):

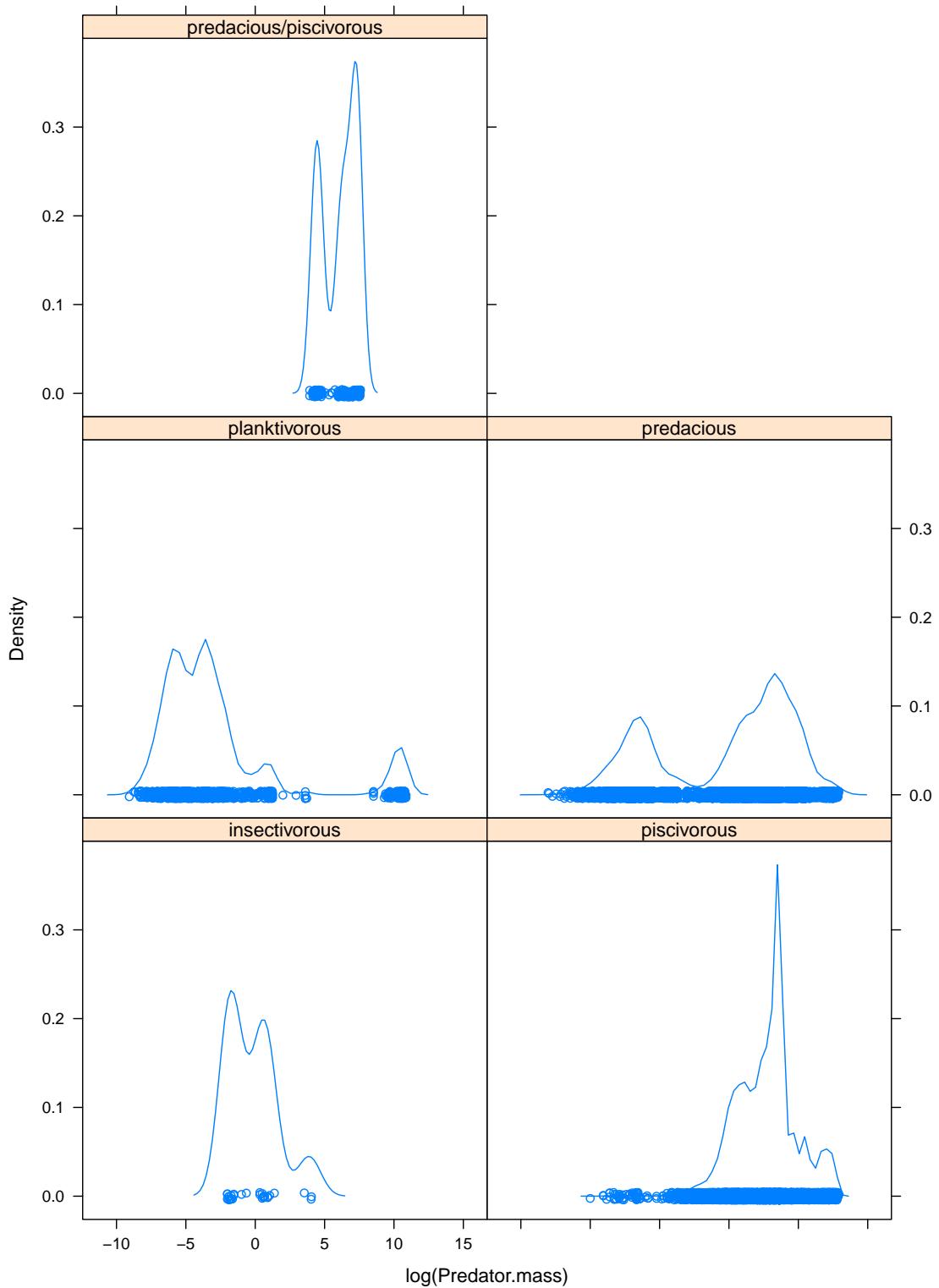


Figure 9.1: A lattice representation of the predator size data

```
> install.packages("ggplot2")
> install.packages("reshape") #A handy additional package
```

### 9.2.1 Basic graphs with qplot

`qplot` stands for quick plot, and is the basic plotting function provided by `ggplot2`. It can be used to quickly produce graphics for exploratory data analysis, and as a base for more complex graphics.

In `ggplot2`, it is necessary to use data frames to store the data. Again we can start plotting the `Predator.mass` vs `Prey.mass`.

```
> require(ggplot2) ## Load the package
Loading required package: ggplot2
> qplot(Prey.mass, Predator.mass, data = MyDF)
```

Again, let's take logarithms and plot:

```
> qplot(log(Prey.mass), log(Predator.mass), data = MyDF)
```

Now, color the points according to the type of feeding interaction:

```
> qplot(log(Prey.mass), log(Predator.mass), data = MyDF,
       colour = Type.of.feeding.interaction)
```

The same as above, but changing the shape:

```
> qplot(log(Prey.mass), log(Predator.mass), data = MyDF,
       shape = Type.of.feeding.interaction)
```

To manually set a color or a shape, you have to use `I()` (meaning “Identity”):

```
> qplot(log(Prey.mass), log(Predator.mass),
       data = MyDF, colour = I("red"))
> qplot(log(Prey.mass), log(Predator.mass),
       data = MyDF, shape= I(3))
```

Because there are so many points, we can make them semi-transparent using `alpha` so that the overlaps can be seen:

```
> qplot(log(Prey.mass), log(Predator.mass), data = MyDF,
       colour = Type.of.feeding.interaction, alpha = I(.5))
```

Now add a smoother to the points.

```
> qplot(log(Prey.mass), log(Predator.mass), data = MyDF,
```

```
geom = c("point", "smooth"))
```

If we want to have a linear regression, we can specify the method `lm`:

```
> qplot(log(Prey.mass), log(Predator.mass), data = MyDF,
       geom = c("point", "smooth"), method = "lm")
```

We can add a smoother for each type of interaction:

```
> qplot(log(Prey.mass), log(Predator.mass), data = MyDF,
       geom = c("point", "smooth"), method = "lm",
       colour = Type.of.feeding.interaction)
```

To extend the lines to the full range, use `fullrange = TRUE`:

```
> qplot(log(Prey.mass), log(Predator.mass), data = MyDF,
       geom = c("point", "smooth"), method = "lm",
       colour = Type.of.feeding.interaction,
       fullrange = TRUE)
```

Now we want to see how the ratio between prey and predator mass changes according to the type of interaction:

```
> qplot(Type.of.feeding.interaction,
       log(Prey.mass/Predator.mass), data = MyDF)
```

Because there are so many points, we can “jitter” them to get a better idea of the spread:

```
> qplot(Type.of.feeding.interaction,
       log(Prey.mass/Predator.mass), data = MyDF,
       geom = "jitter")
```

Or we can draw a boxplot of the data (note the `geom` argument):

```
> qplot(Type.of.feeding.interaction,
       log(Prey.mass/Predator.mass), data = MyDF,
       geom = "boxplot")
```

Now let's draw an histogram of predator-prey mass ratios:

```
> qplot(log(Prey.mass/Predator.mass), data = MyDF,
       geom = "histogram")
```

Color the histogram according to the interaction type:

```
> qplot(log(Prey.mass/Predator.mass), data = MyDF,
       geom = "histogram",
       fill = Type.of.feeding.interaction)
```

You may want to define binwidth (in units of x axis):

```
> qplot(log(Prey.mass/Predator.mass), data = MyDF,
       geom = "histogram",
       fill = Type.of.feeding.interaction,
       binwidth = 1)
```

To make it easier to read, we can plot the smoothed density of the data:

```
> qplot(log(Prey.mass/Predator.mass), data = MyDF,
       geom = "density", fill = Type.of.feeding.interaction)
```

And you can make the densities transparent so that the overlaps are visible:

```
> qplot(log(Prey.mass/Predator.mass), data = MyDF,
       geom = "density", fill = Type.of.feeding.interaction, alpha =
       I(0.5))
```

Or using colour instead of fill draws only the edge of the curve:

```
> qplot(log(Prey.mass/Predator.mass), data = MyDF,
       geom = "density", colour = Type.of.feeding.interaction)
```

Similarly, `geom = "bar"` produces a barplot, `geom = "line"` a series of points joined by a line, etc.

An alternative way of displaying data belonging to different classes is using “faceting”. A simple example:

```
> qplot(log(Prey.mass/Predator.mass),
       facets = Type.of.feeding.interaction ~ .,
       data = MyDF, geom = "density")
```

A more elegant way of drawing logarithmic quantities is to set the axes to be logarithmic:

```
> qplot(Prey.mass, Predator.mass, data = MyDF, log="xy")
```

Let's add a title and labels:

```
> qplot(Prey.mass, Predator.mass, data = MyDF, log="xy",
       main = "Relation between predator and prey mass",
       xlab = "log(Prey mass) (g)",
       ylab = "log(Predator mass) (g)")
```

Adding `+ theme_bw()` makes it suitable for black and white printing.

```
> qplot(Prey.mass, Predator.mass, data = MyDF, log="xy",
       main = "Relation between predator and prey mass",
```

```
xlab = "Prey mass (g)",
ylab = "Predator mass (g)" + theme_bw()
```

Finally, let's save a pdf file of the figure (same approach as we used before):

```
> pdf("../Results/MyFirst-ggplot2-Figure.pdf")
> print(qplot(Prey.mass, Predator.mass, data = MyDF, log="xy",
+   main = "Relation between predator and prey mass",
+   xlab = "log(Prey mass) (g)",
+   ylab = "log(Predator mass) (g)") + theme_bw())
> dev.off()
```

Using `print` ensures that the whole command is kept together and that you can use the command in a script.

Other important options to keep in mind:

<code>xlim</code>	limits for x axis: <code>xlim = c(0,12)</code>
<code>ylim</code>	limits for y axis
<code>log</code>	log transform variable <code>log = "x"</code> , <code>log = "y"</code> , <code>log = "xy"</code>
<code>main</code>	title of the plot <code>main = "My Graph"</code>
<code>xlab</code>	x-axis label
<code>ylab</code>	y-axis label
<code>asp</code>	aspect ratio <code>asp = 2</code> , <code>asp = 0.5</code>
<code>margins</code>	whether or not margins will be displayed

### 9.2.2 Various geom

`geom` Specifies the geometric objects that define the graph type. The `geom` option is expressed as a character vector with one or more entries. `geom` values include “point”, “smooth”, “boxplot”, “line”, “histogram”, “density”, “bar”, and “jitter”. Try the following:

```
# load the package
require(ggplot2)

# load the data
MyDF <- as.data.frame(
  read.csv("../Data/EcolArchives-E089-51-D1.csv"))

# barplot
qplot(Predator.lifestage,
      data = MyDF, geom = "bar")

# boxplot
qplot(Predator.lifestage, log(Prey.mass),
      data = MyDF, geom = "boxplot")

# density
qplot(log(Predator.mass),
      data = MyDF, geom = "density")

# histogram
qplot(log(Predator.mass),
```

```

    data = MyDF, geom = "histogram")

# scatterplot
qplot(log(Predator.mass), log(Prey.mass),
      data = MyDF, geom = "point")

# smooth
qplot(log(Predator.mass), log(Prey.mass),
      data = MyDF, geom = "smooth")

qplot(log(Predator.mass), log(Prey.mass),
      data = MyDF, geom = "smooth", method = "lm")

```

### 9.2.3 Practical 2.2

In this practical, you will write script that draws and saves a pdf file of Fig. 9.2.3, and writes the accompanying regression results to a formatted table in csv. Note that the plots show that the analysis must be subsetted by the `Predator.lifestage` field of the dataset. The guidelines are:

- ★ Write a script file called `PP_Regress.R` and save it in the `Code` directory — sourcing or running this script should result in one pdf file containing (Fig. 9.2.3) being saved in the `Results` directory (Hint: Use the `print()` command to write to the pdf).
- ★ In addition, the script should calculate the regression results corresponding to the lines fitted in Fig. 9.2.3 and save it to a csv delimited table called (`PP_Regress_Results.csv`), in the `Results` directory. (Hint: you will have to initialize a new dataframe in the script to first store the calculations and then `write.csv()` or `write.table()` it.)  
All that you are being asked for here is results of an analysis of Linear regression on subsets of the data corresponding to available Feeding Type  $\times$  Predator life Stage combination — not a multivariate linear model with these two as separate covariates!
- ★ The regression results should include the following with appropriate headers (e.g., slope, intercept, etc, in each Feeding type  $\times$  life stage category): regression slope, regression intercept,  $R^2$ , F-statistic value, and p-value of the overall regression (Hint: Review Practical 8 of the Stats week).
- ★ The script should be self-sufficient and not need any external inputs — it should import the above predator-prey dataset from the appropriate directory, and save the graphic plots to the appropriate directory (Hint: use relative paths). I should be able to `source` it without errors.

**Extra Credit:** Do the same as above, but the analysis this time should be separate by the dataset's `Location` field. This is a substantial extra effort, so I do think it deserves substantial extra credit!

### 9.2.4 Advanced plotting: `ggplot`

The command `qplot` allows you to use only a single dataset and a single set of “aesthetics” (x, y, etc.). To make full use of `ggplot2`, we need to use the command `ggplot`. We need:

- The data to be plotted, in a data frame;
- Aesthetics mappings, specifying which variables we want to plot, and how;
- The `geom`, defining how to draw the data;

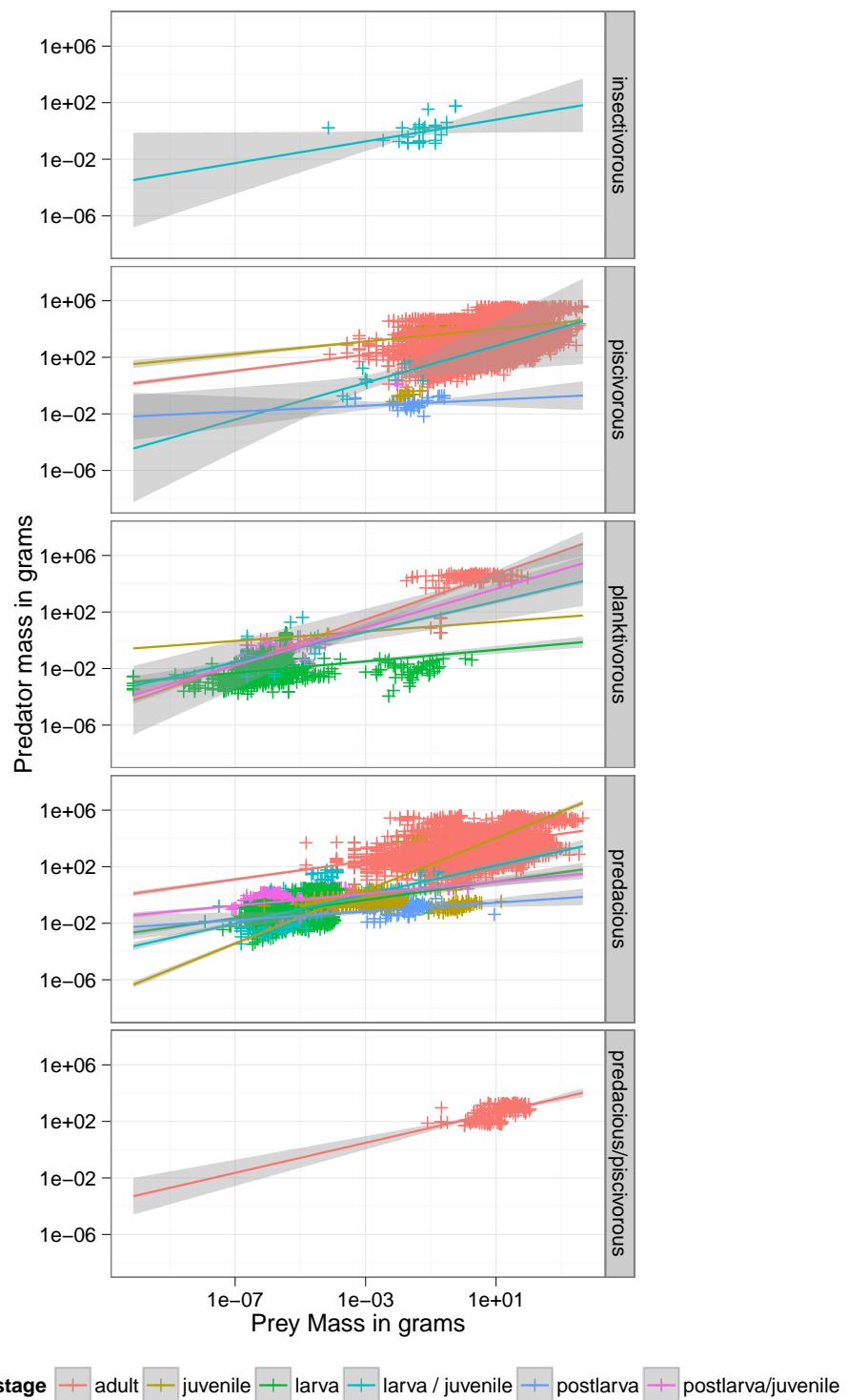


Figure 9.2: Write a script that generates this figure.

- (Optionally) some `stat` that transform the data or perform statistics using the data.

To start a graph, we can specify the data and the aesthetics:

```
> p <- ggplot(MyDF, aes(x = log(Predator.mass),
                           y = log(Prey.mass),
                           colour = Type.of.feeding.interaction ))
```

Now try to plot the graph:

```
> p
Error: No layers in plot
```

In fact, we have to specify a geometry in order to see the graph:

```
> p + geom_point()
```

We can use the “plus” sign to concatenate different commands:

```
> p <- ggplot(MyDF, aes(x = log(Predator.mass),
                           y = log(Prey.mass),
                           colour = Type.of.feeding.interaction ))
> q <- p + geom_point(size=I(2), shape=I(10)) + theme_bw()
> q
```

Let's remove the legend:

```
> q + opts(legend.position = "none")
```

### 9.2.5 Case study 1: plotting a matrix

In this section we will plot a matrix of random values taken from a normal distribution  $\mathcal{U}[0, 1]$ . Our goal is to produce the plot in Figure 9.3. Because we want to plot a matrix, and `ggplot2` accepts only dataframes, we use the package `reshape` that can “melt” a matrix into a dataframe:

```
require(ggplot2)
require(reshape)

GenerateMatrix <- function(N) {
  M <- matrix(runif(N * N), N, N)
  return(M)
}

> M <- GenerateMatrix(10)

> M[1:3, 1:3]
      [,1]      [,2]      [,3]
[1,] 0.2700254 0.8686728 0.7365857
[2,] 0.1744879 0.8488169 0.4165879
[3,] 0.3980783 0.7727821 0.4271121
```

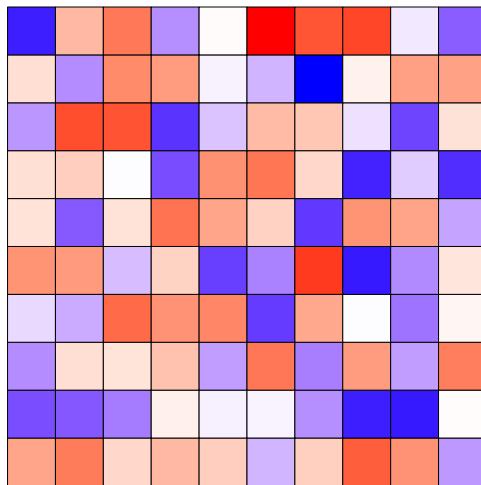


Figure 9.3: Random matrix with values sampled from uniform distribution.

```

> Melt <- melt(M)

> Melt[1:4,]
   X1 X2      value
1  1  1 0.2700254
2  2  1 0.1744879
3  3  1 0.3980783
4  4  1 0.3196671

> ggplot(Melt, aes(X1, X2, fill = value)) + geom_tile()

# adding a black line dividing cells
> p <- ggplot(Melt, aes(X1, X2, fill = value))
> p <- p + geom_tile(colour = "black")

# removing the legend
> q <- p + opts(legend.position = "none")

# removing all the rest
> q <- p + opts(legend.position = "none",
+                 panel.background = theme_blank(),
+                 axis.ticks = theme_blank(),
+                 panel.grid.major=theme_blank(),
+                 panel.grid.minor=theme_blank(),
+                 axis.text.x = theme_blank(),
+                 axis.title.x=theme_blank(),
+                 axis.text.y = theme_blank(),
+                 axis.title.y=theme_blank())

# exploring the colors
> q + scale_fill_continuous(low = "yellow",
+                            high = "darkgreen")
> q + scale_fill_gradient2()
> q + scale_fill_gradientn(colours = grey.colors(10))
> q + scale_fill_gradientn(colours = rainbow(10))
> q + scale_fill_gradientn(colours =

```

```
c("red", "white", "blue"))
```

### 9.2.6 Case study 2: plotting two dataframes

According to Girko's circular law, the eigenvalues of a matrix  $M$  of size  $N \times N$  are approximately contained in a circle in the complex plane with radius  $\sqrt{N}$ . We are going to draw a simulation displaying this result (Figure 9.4).

```
require(ggplot2)

# function that returns an ellipse
build_ellipse <- function(hradius, vradius) {
  npoints = 250
  a <- seq(0, 2 * pi, length = npoints + 1)
  x <- hradius * cos(a)
  y <- vradius * sin(a)
  return(data.frame(x = x, y = y))
}

# Size of the matrix
N <- 250
# Build the matrix
M <- matrix(rnorm(N * N), N, N)
# Find the eigenvalues
eigvals <- eigen(M)$values
# Build a dataframe
eigDF <- data.frame("Real" = Re(eigvals),
                     "Imaginary" = Im(eigvals))

# The radius of the circle is sqrt(N)
my_radius <- sqrt(N)
# Ellipse dataframe
ellDF <- build_ellipse(my_radius, my_radius)
# rename the columns
names(ellDF) <- c("Real", "Imaginary")

# Now the plotting:
# plot the eigenvalues
p <- ggplot(eigDF, aes(x = Real, y = Imaginary))
p <- p +
  geom_point(shape = I(3)) +
  opts(legend.position = "none")

# now add the vertical and horizontal line
p <- p + geom_hline(aes(intercept = 0))
p <- p + geom_vline(aes(intercept = 0))

# finally, add the ellipse
p <- p + geom_polygon(data = ellDF,
                       aes(x = Real,
                           y = Imaginary,
                           alpha = 1/20,
                           fill = "red"))

pdf("Girko.pdf")
print(p)
dev.off()
```

### 9.2.7 Case study 3: annotating the plot

In the plot in Figure 9.5, we use the geometry “text” to annotate the plot.

```

require(ggplot2)

filename <- "Results.txt"
a <- read.table(filename, header = TRUE)
# here's how the data looks like
print(a[1:3,])
print(a[90:95,])

# append a col of zeros
a$ymin <- rep(0, dim(a)[1])

# print the first linerange
p <- ggplot(a)
p <- p + geom_linerange(data = a, aes(
  x = x,
  ymin = ymin,
  ymax = y1,
  size = (0.5)
),
colour = "#E69F00",
alpha = 1/2, show_guide = FALSE)

# print the second linerange
p <- p + geom_linerange(data = a, aes(
  x = x,
  ymin = ymin,
  ymax = y2,
  size = (0.5)
),
colour = "#56B4E9",
alpha = 1/2, show_guide = FALSE)

# print the third linerange
p <- p + geom_linerange(data = a, aes(
  x = x,
  ymin = ymin,
  ymax = y3,
  size = (0.5)
),
colour = "#D55E00",
alpha = 1/2, show_guide = FALSE)

# annotate the plot with labels
p <- p + geom_text(data = a,
  aes(x = x, y = -500, label = Label))

# now set the axis labels,
# remove the legend, prepare for bw printing
p <- p + scale_x_continuous("My x axis",
  breaks = seq(3, 5, by = 0.05)
)
scale_y_continuous("My y axis") + theme_bw() +
opts(legend.position = "none")

```

```
# Finally, print in a pdf
pdf("MyBars.pdf", width = 12, height = 6)
print(p)
dev.off()
```

### 9.2.8 Case study 4: mathematical display

In Figure 9.6, you can see the mathematical annotation of the axis and on the plot.

```
require(ggplot2)

# create an "ideal" linear regression data!
x <- seq(0, 100, by = 0.1)
y <- -4. + 0.25 * x +
  rnorm(length(x), mean = 0., sd = 2.5)

# now a dataframe
my_data <- data.frame(x = x, y = y)

# perform a linear regression
my_lm <- summary(lm(y ~ x, data = my_data))

# plot the data
p <- ggplot(my_data, aes(x = x, y = y,
                           colour = abs(my_lm$residual)))
  ) +
  geom_point() +
  scale_colour_gradient(low = "black", high = "red") +
  opts(legend.position = "none") +
  scale_x_continuous(
    expression(alpha^2 * pi / beta * sqrt(Theta)))

# add the regression line
p <- p + geom_abline(
  intercept = my_lm$coefficients[1][1],
  slope = my_lm$coefficients[2][1],
  colour = "red")
# throw some math on the plot
p <- p + geom_text(aes(x = 60, y = 0,
                        label = "sqrt(alpha) * 2* pi"),
                    parse = TRUE, size = 6,
                    colour = "blue")

# print in a pdf
pdf("MyLinReg.pdf")
print(p)
dev.off()
```

## 9.3 Readings

- The classic Tufte [www.edwardtufte.com/tufte/books\\_vdq](http://www.edwardtufte.com/tufte/books_vdq) (btw, check out what Tufte thinks of PowerPoint!)

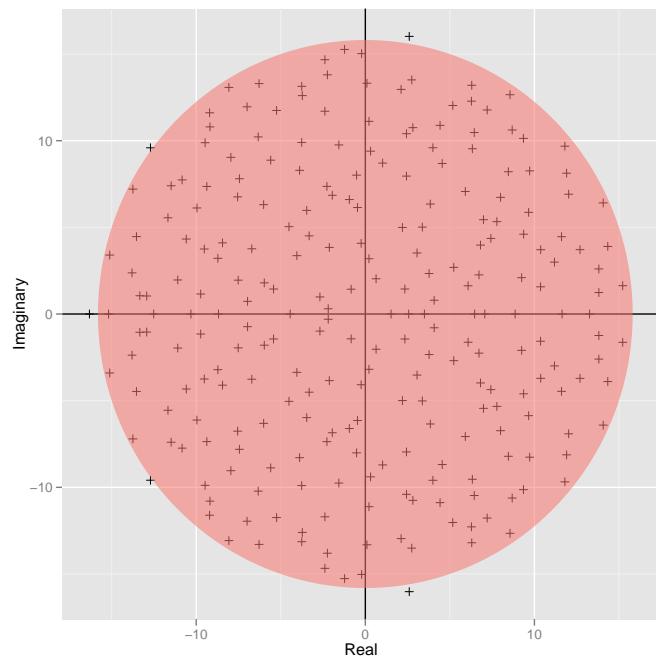


Figure 9.4: Girko's circular law.

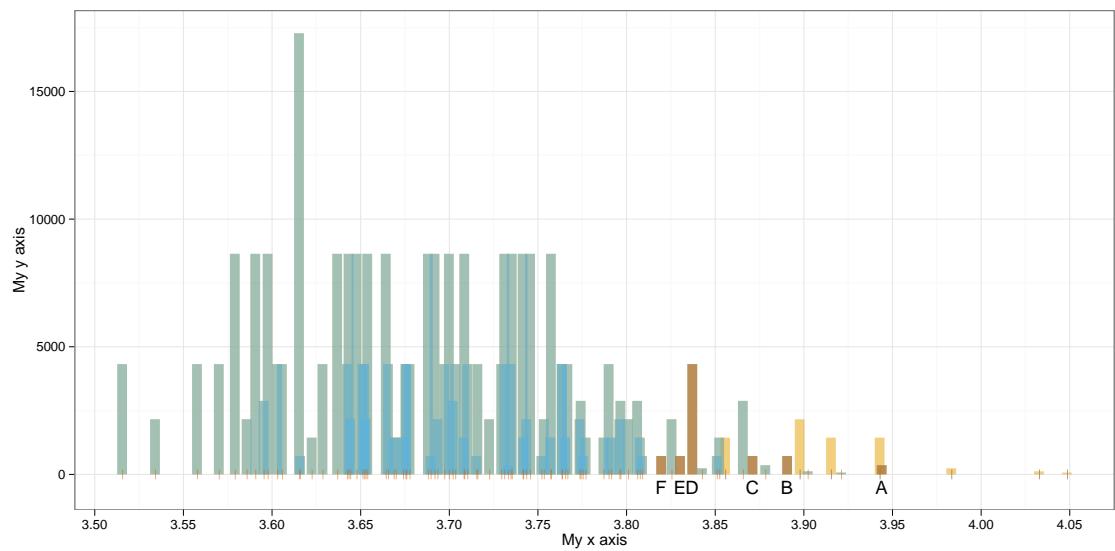


Figure 9.5: Overlay of three lineranges and a text geometry.

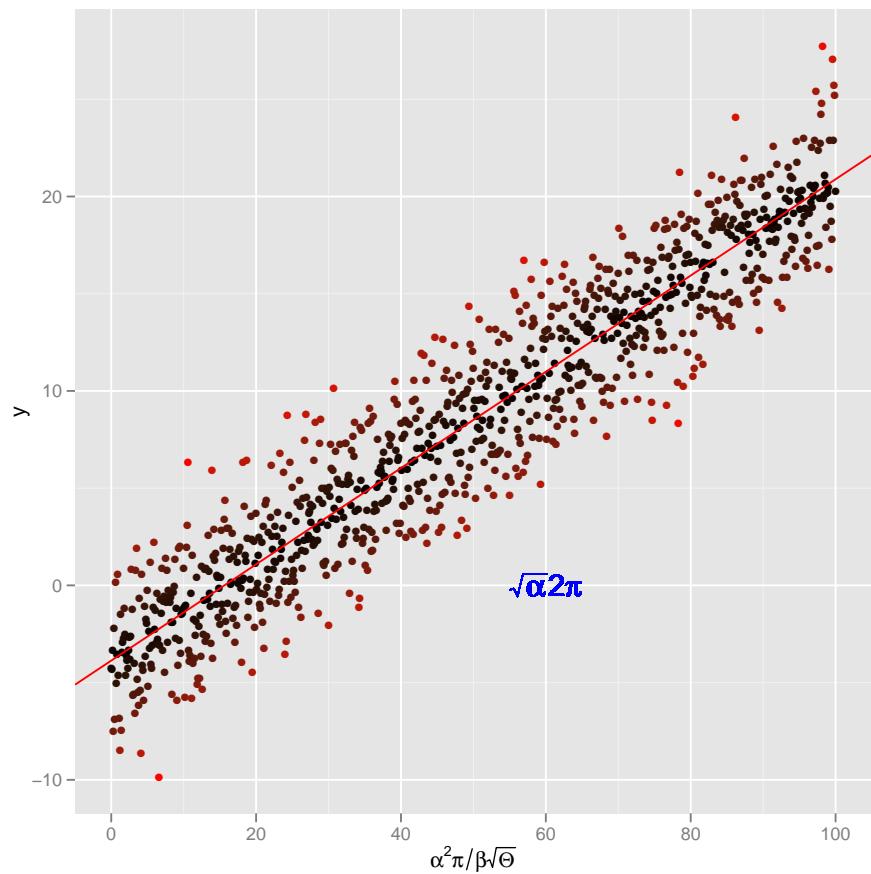


Figure 9.6: Linear regression with colors expressing residuals and mathematical annotations.

Available in the Central Library, I have also added extracts and a related book in pdf (on Blackboard)

- Rolandi et al. “A Brief Guide to Designing Effective Figures for the Scientific Paper”, doi:10.1002/adma.201102518 (on Blackboard)
- Lauren et al. “Graphs, Tables, and Figures in Scientific Publications: The Good, the Bad, and How Not to Be the Latter”, doi:10.1016/j.jhsa.2011.12.041 (on Blackboard)
- Effective scientific illustrations: [www.labtimes.org/labtimes/issues/lt2008/lt05/lt\\_2008\\_05\\_52\\_53.pdf](http://www.labtimes.org/labtimes/issues/lt2008/lt05/lt_2008_05_52_53.pdf) (on Blackboard)
- <https://web.archive.org/web/20120310121708/http://addictedor.free.fr/graphiques/thumbs.php>
- Make xkcd style graphs in R: <http://xkcd.r-forge.r-project.org/>

# Chapter 10

## Advanced topics in R

### 10.1 Vectorization

R is very slow at running cycles (`for` and `while` loops). This is because R is a “nimble” language: at execution time R does not know what you’re going to perform until it “reads” the code to perform. Compiled languages such as C, know exactly what the flow of the program is, as the code is compiled before execution. As a metaphor, C is a musician playing a score she has seen before – optimizing each passage, while R is playing it “a prima vista” (i.e., at first sight).

Hence, in R you should try to avoid loops like the plague. In practical terms, sometimes it is much easier to throw in a `for` loop, and then optimize the code to avoid the loop if the running time is not satisfactory. R has several functions that can operate on entire vectors and matrices.

- \* For example, type (save in `Code`) in `Vectorize1.R` and run it (it sums all elements of a matrix):

```
M <- matrix(runif(1000000),1000,1000)

SumAllElements <- function(M) {
  Dimensions <- dim(M)
  Tot <- 0
  for (i in 1:Dimensions[1]){
    for (j in 1:Dimensions[2]){
      Tot <- Tot + M[i,j]
    }
  }
  return (Tot)
}

## This on my computer takes about 1 sec
print(system.time(SumAllElements(M)))
## While this takes about 0.01 sec
print(system.time(sum(M)))
```

Both approaches are correct, and will give you the right answer. However, one is 100 times faster than the other!

Fortunately, R offers several ways of avoiding loops. Here are the main ones:

```

## apply:
# applying the same function to rows/columns of a matrix

## Build a random matrix
M <- matrix(rnorm(100), 10, 10)
## Take the mean of each row
RowMeans <- apply(M, 1, mean)
print (RowMeans)
## Now the variance
RowVars <- apply(M, 1, var)
print (RowVars)
## By column
ColMeans <- apply(M, 2, mean)
print (ColMeans)

## You can use it to define your own functions
## (What does this function do?)
SomeOperation <- function(v) {
  if (sum(v) > 0) {
    return (v * 100)
  }
  return (v)
}
print (apply(M, 1, SomeOperation))

## by:
## apply the function to a dataframe, using some factor
## to define the subsets

## import some data
attach(iris)
print (iris)

## use colMeans (as it is better for dataframes)
by(iris[,1:2], iris$Species, colMeans)
by(iris[,1:2], iris$Petal.Width, colMeans)

## There are many other methods: lapply, sapply, eapply, etc.
## Each is best for a given data type (lapply -> lists)

## replicate:
## this is quite useful to avoid a loop for function that typically
## involve random number generation
print(replicate(10, runif(5)))

```

## 10.2 Breaking out of loops

Often it is useful (or necessary) to break out of a loop when some condition is met. Use `break` (like in pretty much any other programming language, like python) in situations when you cannot set a target number of iterations, as you would with a `while` loop (Chapter 1). Try this (type into `break.R` and save in Code):

```

i <- 0 #Initialize i
while(i < Inf) {
  if (i == 20) {

```

```

        break } # Break out of the while loop!
    else {
        cat("i equals " , i , " \n")
        i <- i + 1 # Update i
    }
}

```

## 10.3 Using next

You can also skip to next iteration of a loop. Both `next` and `break` can be used within other loops (`while`, `for`). Try this (type into `next.R` and save in `Code` (what does this script do?)):

```

for (i in 1:10) {
  if ((i %% 2) == 0)
    next # pass to next iteration of loop
  print(i)
}

```

**Reminder:** Indent your code! Indentation helps you see the flow of the logic, rather than flattened version, which is hard for you and everybody else to read. I recommend using the `tab` key to indent.

## 10.4 “Catching” errors

Often, you don’t know if a simulation or a R function will work on a particular data or variable, or a value of a variable (can happen in many stats functions). Rather than having R throw you out of the code, you would rather catch the error and keep going. This can be done using `try`. Type the following into `try.R` and save in `Code` (what does this script do?):

```

## run a simulation that involves sampling from a population

x <- stats::rnorm(50)
doit <- function(x){
  x <- sample(x, replace = TRUE)
  if(length(unique(x)) > 30) {
    mean(x)
  } else {
    stop("too few unique points")
  }
}

## Try using "try" with vectorization:
result <- lapply(1:100, function(i) try(doit(x), FALSE))

## Or using a for loop:
result <- vector("list", 100) #Preallocate/Initialize
for(i in 1:100) {
  result[[i]] <- try(doit(x), FALSE)
}

```

Note the functions `sample` and `stop` in the above script. Also check out `tryCatch`.

## 10.5 Generating Random Numbers

Computers don't really generate mathematically random numbers, but instead a sequence of numbers that are close to random: “pseudo-random numbers”. They are generated based on some iterative formula:

$$x_{new} = f(x_{old}) \mod N$$

where modulo operation provides the “remainder” division.

To generate the first random number, you need a **seed**. Setting the seed allows you to reliably generate the same sequence of numbers, which can be useful when debugging programs (next section).

R has many routines for generating random samples from various probability distributions — we have already used `rnorm()`, `runif()`. Try this:

```
> set.seed(1234567)
> rnorm(1)
0.1567038
```

What happened?! If this were truly a random number, how would everybody get the same answer? Now try `rnorm(10)` and compare the results with your neighbour. Thus “random” numbers generated in R and in any other software are in fact “deterministic”, but from a very complex formula that yields numbers with properties like random numbers.

Effectively, `rnorm` has an enormous list that it cycles through. The random seed starts the process, i.e., indicates where in the list to start. This is usually taken from the clock when you start R.

But why bother with this? Well, for debugging (next section). Bugs in code can be hard to find — harder still if you are generating random numbers, so repeat runs of your code may or may not all trigger the same behaviour. You can set the seed once at the beginning of the code — ensuring repeatability, retaining (pseudo) randomness. Once debugged, if you want, you can remove the set seed line.

## 10.6 Debugging

Indeed, as most of you must have already experienced by now, there can be frustrating, puzzling bugs in programs that lead to mysterious errors. Often, the error and warning messages you get are un-understandable, especially in R! Some useful debugging functions in R:

- Warnings vs Errors; converting warnings to errors: `stopifnot()` — a bit like `try`
- What to do when you get an error: `traceback()`
- Simple `print` commands in the right places can be useful for testing (but not strongly recommended)
- Use of `browser()` at key points in code — my favourite option (also look up `recover()`)
- `debug(fn)`, `undebug(fn)` : More technical approach to debugging — explore them

Let's look at an example using `browser()`. `browser()` is handy because it will allow you to "single-step" through your code. Place it within your function at the point you want to examine (e.g.) local variables.

Here's an example usage of `browser()` (type in `browse.R` and save in Code):

```
Exponential <- function(N0 = 1, r = 1, generations = 10) {
  # Runs a simulation of exponential growth
  # Returns a vector of length generations

  N <- rep(NA, generations)      # Creates a vector of NA

  N[1] <- N0
  for (t in 2:generations) {
    N[t] <- N[t-1] * exp(r)
    # browser()
  }
  return (N)
}
plot(Exponential(), type="l", main="Exponential growth")
```

Now, within the browser, you can enter expressions as normal, or you can use a few particularly useful debug commands:

- n: single-step
- c: exit browser and continue
- Q: exit browser and abort, return to top-level.

## 10.7 Practical

*Non-CMEE students can ignore this one!*

The Ricker model is a classic discrete population model which was introduced in 1954 by Ricker to model recruitment of stock in fisheries. It gives the expected number (or density)  $N_{t+1}$  of individuals in generation  $t+1$  as a function of the number of individuals in the previous generation  $t$ :

$$N_{t+1} = N_t e^{r(1 - \frac{N_t}{K})} \quad (10.1)$$

Here  $r$  is intrinsic growth rate and  $k$  as the carrying capacity of the environment. Try this script that runs it:

```
Ricker <- function(N0=1, r=1, K=10, generations=50)
{
  # Runs a simulation of the ricker model
  # Returns a vector of length generations

  N <- rep(NA, generations)      # Creates a vector of NA

  N[1] <- N0
  for (t in 2:generations)
  {
    N[t] <- N[t-1] * exp(r*(1.0-(N[t-1]/K)))
  }
}
```

```

    return (N)
}
plot(Ricker(generations=10), type="l")

```

Now open and run the script `Vectorize2.R` (available on the bitbucket Git repository). This is the stochastic Ricker model (compare with the above script to see where the stochasticity (random error) enters). Now modify the script to complete the exercise given. CMEEs, As always, bring your functional code and data under version control!

## 10.8 Launching/Running R in batch mode

Often, you may want to run the final analysis without opening R in interactive mode. In Mac or linux, you can do so by typing:

```
R CMD BATCH MyCode.R MyResults.Rout
```

This will create an `MyResults.Rout` file containing all the output. On Microsoft Windows, it's more complicated — change the path to `R.exe` and output file as needed:

```
"C:\Program Files\R\R-3.1.1\bin\R.exe" CMD BATCH -vanilla -slave
"C:\PathToMyResults\Results\MyCode.R"
```

## 10.9 Accessing databases using R

R can be used to link seamlessly to on-line databases such as PubMed and GenBank. Such computational Biology datasets are often quite large, and it makes sense to access their data by querying the databases instead of manually downloading. There are many R packages that provide an interface to databases (SQLite, MySQL, Oracle, etc). Check out R packages DBI (<http://cran.r-project.org/web/packages/DBI/index.html>) and RMySQL ([cran.r-project.org/web/packages/RMySQL/index.html](http://cran.r-project.org/web/packages/RMySQL/index.html)).

## 10.10 Building your own R packages

You can packaging up your code, data sets and documentation to make a *bona fide* R package. You may wish to do this for particularly large projects that you think will be useful for others. Read *Writing R Extensions* ([cran.r-project.org/doc/manuals/r-release/R-exts.html](http://cran.r-project.org/doc/manuals/r-release/R-exts.html)) manual and see `package.skeleton` to get started. The R tool set EcoDataTools (<https://github.com/DomBennett/EcoDataTools>) and the package `cheddar` were written by Silwood Grad Students!

## 10.11 Sweave

Sweave is a tool that allows you to write your Dissertation or Report such that it can be updated automatically if data or R analysis change. Instead of inserting a prefabricated graph or table into the report, the master document contains the R code necessary to obtain it. When run through R,

all data analysis output (tables, graphs, etc.) is created on the fly and inserted into a final latex document. The report can be automatically updated if data or analysis change, which allows for truly reproducible research. Check out the webpage: <http://www.statistik.lmu.de/~leisch/Sweave>

## 10.12 Practical (CMEE Extra Credit)

Autocorrelation in weather

- \* Make a new script named `TAutoCorr.R`, and save in `Code` directory
- \* At the start of the script, load and examine and plot `KeyWestAnnualMeanTemperature.Rdata`, using `load()` — This is the temperature in Key West, Florida for the 20th century
- \* The question this script will help answer is: Is the temperature one year significantly correlated with the next year?
- \* You will need the `sample` function that we used today — read the help file for `sample` and experiment with it
- \* Now,
  1. Compute the appropriate correlation coefficient and store it (look at the help file for `cor()`)
  2. Then repeat this 10000 times
    - Randomly permute the time series
    - Compute the same correlation coefficient
    - Store it
  3. Then see what fraction of the coefficients from step 2 were greater than that from step 1
- \* How do you interpret these results? Why? Present your interpretation in a pdf document.

## 10.13 Practical (CMEE Extra Credit)

Your project may not really need GIS, but you would still like to do some mapping. Yo you can do it in R using the `maps` package. In this practical, you will map the Global Population Dynamics Database (<https://www.imperial.ac.uk/cpb/gpdd/gpdd.aspx>) (GPDD). This is a freely available database that was developed at Silwood).

If any of you are interested in doing a project around this database, please contact David Orme or Samraat Pawar! It is a goldmine of as yet under-utilized information. Note that the Living Planet Index (<http://livingplanetindex.org/home/index>) is based upon these data.

- \* Use `load()` from `GPDDFiltered.RData` that is available on the Git repository — have a look at the database field headers and contents.
- \* What you need is latitude and logitude information for a bunch of species for which population time series are available in the GPDD
- \* Now use `install.packages()` to install the package `maps`, as you did with `ggplot2` — hopefully without any problems!
- \* Now create a script (saved under a sensible name in a sensible location — hint hint!) that:
  1. Loads the `maps` package
  2. Loads the GPDD data

3. Creates a world map (use the map function, read its help file, also google examples using maps)
4. Superimposes on the map all the locations from which we have data in the GPDD dataframe
5. Compare your map with a fellow student to check

*Based on this map, what biases might you expect in any analysis based on the data represented?*

## 10.14 R Module Wrap up

### 10.14.1 Some comments and suggestions

Thanks for enduring through the week! Learning to program in R or any other language, especially if it's your first-ever effort to learn programming, demands perseverance. Most of you have shown admirable quantities of this necessary quality. Keep going! I believe most of you have climbed a significant part of a steep learning curve. Here are some things to keep in mind:

- There are many R nerds at Silwood that you can talk to — They walk among us!
- There is a Silwood R list that you can subscribe to: <https://mailman.ic.ac.uk/mailman/listinfo/silwood-r>
- However, post questions only as a last resort! Google it first, and even before that, make sure you revise this week's (and stats week's) work.
- Solutions to this weeks Pracs/Exercises will become available by Thursday (13th Nov) next week.

### 10.14.2 CMEE weekly coursework submission

Test and bring under version control: TreeHeight.R, control.R, PP\_Lattice.R, PP\_Regress.R (where's the extra credit?), TAtoCorr.R, Vectorize1.R, Ricker.R, Vectorize2.R, break.R, next.R, try.R, browse.R, TAtoCorr.R (Extra Credit).

Git commit and push by: **Wednesday, 12 November 2014**

## 10.15 Readings

- See **An introduction to the Interactive Debugging Tools in R, Roger D Peng** for detailed usage. <http://www.biostat.jhsph.edu/~rpeng/docs/R-debug-tools.pdf>
- Friedrich Leisch. (2002) Sweave: Dynamic generation of statistical reports using literate data analysis. *Proceedings in Computational Statistics*, pages 575-580. Physica Verlag, Heidelberg, 2002. <http://www.statistik.lmu.de/~leisch/Sweave>
- Remember, R packages come with pdf guides/documentation!

