

# Merely Useful

*Dhavid Aruliah, Madeleine Bonsma-Fisher, Jonathan Dursi, Kate Hertweck, Katy Huff, Damien Irving, Luke Johnston, Christina Koch, Sara Mahallati, Brandeis Marshall, Joel Ostblom, Matt Turk, Elizabeth Wickes, Charlotte Wickham, and Greg Wilson*

*2019-09-12*



# Contents

<b>1</b>	<b>Overview</b>	<b>15</b>
1.1	Why isn't all of this normal already? . . . . .	15
1.2	Acknowledgments . . . . .	16
<b>2</b>	<b>Novice Goals</b>	<b>17</b>
2.1	Personas . . . . .	17
2.2	Getting started . . . . .	19
2.3	Data manipulation . . . . .	20
2.4	Plotting . . . . .	20
2.5	Development . . . . .	21
2.6	Data analysis . . . . .	21
2.7	Version control . . . . .	22
2.8	Publishing . . . . .	23
2.9	Reproducibility . . . . .	23
2.10	Collaboration . . . . .	24
<b>I</b>	<b>Novice R Material</b>	<b>25</b>
<b>3</b>	<b>Introduction</b>	<b>27</b>
3.1	Who are these lessons for? . . . . .	27
3.2	What will these lessons teach you? . . . . .	28
3.3	What examples will we use? . . . . .	29
3.4	What's the big picture? (#novice-r-intro-bigpicture) . . . . .	29
<b>4</b>	<b>Practice</b>	<b>31</b>
4.1	Working with a single tidy table . . . . .	31
4.2	Working with grouped data . . . . .	34
4.3	Creating charts . . . . .	39
<b>5</b>	<b>Reproducibility</b>	<b>53</b>
5.1	Questions {r-reproducibility-questions} . . . . .	53
5.2	Objectives . . . . .	53
5.3	Introduction . . . . .	53

5.4	Project organization . . . . .	54
5.5	Reusability . . . . .	57
5.6	Readability . . . . .	59
5.7	Integrating text, code, and results . . . . .	64
5.8	Key Points . . . . .	79
5.9	Additional learning resources and material . . . . .	79
<b>6</b>	<b>Data Manipulation</b>	<b>81</b>
6.1	Questions . . . . .	81
6.2	Motivation . . . . .	81
6.3	Exploring data in the console . . . . .	84
6.4	How can I select subsets of my data? . . . . .	86
6.5	How can I calculate new values? . . . . .	96
6.6	How can I tell what's gone wrong in my programs? . . . . .	100
6.7	How can I operate on subsets of my data? . . . . .	107
6.8	How can I read my own tabular data into R? . . . . .	112
6.9	How can I save my results? . . . . .	114
<b>7</b>	<b>Publishing</b>	<b>119</b>
7.1	Questions . . . . .	119
7.2	Why should I share my work on the internet? . . . . .	119
7.3	What does it take to get a webpage online? . . . . .	120
7.4	How do I get my work on the web? . . . . .	120
7.5	How do I link to other pages, files or images? . . . . .	126
7.6	Exercise: Add a website to an existing project . . . . .	129
<b>II</b>	<b>Novice Python Material</b>	<b>131</b>
<b>8</b>	<b>Introduction</b>	<b>133</b>
8.1	Who are these lessons for? . . . . .	133
8.2	What does “done” look like? . . . . .	133
8.3	What will we use as running examples? . . . . .	133
<b>9</b>	<b>Development</b>	<b>135</b>
9.1	Questions . . . . .	135
9.2	Objectives . . . . .	135
9.3	Functions . . . . .	135
9.4	How to make programs indicate that something has gone wrong? . . . . .	140
9.5	Packages . . . . .	143
9.6	How to get help online . . . . .	147
9.7	Key Points . . . . .	147
9.8	Exercises . . . . .	148

<b>III RSE Material</b>	<b>149</b>
<b>10 RSE Introduction</b>	<b>151</b>
10.1 Who are these lessons for? . . . . .	152
10.2 What does “done” look like? . . . . .	152
10.3 What will this course accomplish? . . . . .	153
10.4 What will we use as running examples? . . . . .	154
<b>11 Bash Shell</b>	<b>157</b>
11.1 Questions . . . . .	157
11.2 Objectives . . . . .	157
11.3 Introduction . . . . .	159
11.4 Navigating Files and Directories . . . . .	161
11.5 Working With Files and Directories . . . . .	171
11.6 Pipes and Filters . . . . .	177
11.7 Loops . . . . .	181
11.8 Shell Scripts . . . . .	186
11.9 Finding Things . . . . .	190
11.10 Summary . . . . .	197
11.11 Exercises . . . . .	197
11.12 Key Points . . . . .	221
	<b>225</b>
<b>12 A Branching Workflow</b>	<b>227</b>
12.1 Questions . . . . .	227
12.2 Objectives . . . . .	227
12.3 Introduction . . . . .	227
12.4 How can I use branches to manage development of new features? . . . . .	228
12.5 How can I switch between branches when work is only partly done? . . . . .	229
12.6 How can I keep my project’s history clean when working on many branches? . . . . .	230
12.7 How can I make it easy for people to review my work before I merge it? . . . . .	231
12.8 What <i>is</i> a feature? . . . . .	232
12.9 How can I label specific versions of my work? . . . . .	233
12.10 Summary . . . . .	234
12.11 Exercises . . . . .	234
12.12 Key Points . . . . .	234
<b>13 Configuring Software</b>	<b>235</b>
13.1 Questions . . . . .	235
13.2 Objectives . . . . .	235
13.3 Introduction . . . . .	235
13.4 How can I handle command-line flags consistently? . . . . .	236
13.5 What do I do when I run out of memorable single-letter flags? . . . . .	238

13.6	How can I manage configuration files consistently?	238
13.7	How can I implement overlay configuration?	240
13.8	How can I find configuration files?	241
13.9	How can I keep a record of the actual configuration that produced particular results?	242
13.10	Summary	243
13.11	Exercises	243
13.12	Key Points	244
<b>14</b>	<b>Automating Analyses</b>	<b>245</b>
14.1	Questions	245
14.2	Objectives	245
14.3	Introduction	245
14.4	How can I update a file when its prerequisites change?	247
14.5	How can I tell Make where to find rules?	249
14.6	How can I update multiple files when their prerequisites change?	249
14.7	How can I get rid of temporary files that I don't need?	250
14.8	How can I make a target depend on several prerequisites?	251
14.9	How can I reduce the amount of typing I have to do?	251
14.10	How can I make one update depend on another?	252
14.11	How can I abbreviate my update rules?	253
14.12	How can I write one general rule to update many files in the same way?	254
14.13	How can I define sets of files automatically?	255
14.14	How can I document my workflow?	257
14.15	Summary	259
14.16	Exercises	261
14.17	Key Points	261
<b>15</b>	<b>Unit Testing</b>	<b>263</b>
15.1	Questions	263
15.2	Objectives	263
15.3	Introduction	263
15.4	What are realistic goals for testing?	264
15.5	What does a systematic software testing framework look like?	265
15.6	How can I manage tests systematically?	266
15.7	How can I tell if my software failed as it was supposed to?	267
15.8	How can I test software that includes randomness?	269
15.9	How can I test software that does I/O?	270
15.10	How can I tell which parts of my software have (not) been tested?	272
15.11	Summary	273
15.12	Exercises	274
15.13	Key Points	274
<b>16</b>	<b>Verification</b>	<b>275</b>
16.1	Questions	275

16.2 Objectives . . . . .	275
16.3 Introduction . . . . .	275
16.4 What is the difference between testing in software engineering and in data analysis? . . . . .	276
16.5 Why should I be cautious when using floating-point numbers? . .	277
16.6 How can I express how close one number is to another? . . . .	280
16.7 How should I write tests that involved floating-point values? . .	281
16.8 How can I test plots and other graphical results? . . . . .	282
16.9 How can I test the steps in a data analysis pipeline during devel- opment? . . . . .	283
16.10 How can I check the steps in a data analysis pipeline in production?	284
16.11 How can I infer and check properties of my data? . . . . .	286
16.12 Summary . . . . .	287
16.13 Exercises . . . . .	288
16.14 Key Points . . . . .	288
<b>17 Project Structure</b>	<b>289</b>
17.1 Questions . . . . .	289
17.2 Objectives . . . . .	289
17.3 Introduction . . . . .	289
17.4 What are Noble's Rules? . . . . .	289
17.5 How should files and sub-directories be named? . . . . .	290
17.6 How should I manage a mix of compiled programs and scripts? .	292
17.7 Should I separate documentation from manuscripts? . . . . .	292
17.8 How should I handle data can't be stored in version control? . .	293
17.9 What other files should every project contain? . . . . .	293
17.10 What <i>is</i> a project? . . . . .	294
17.11 Summary . . . . .	294
17.12 Exercises . . . . .	295
17.13 Key Points . . . . .	295
<b>18 Including Everyone</b>	<b>297</b>
18.1 Questions . . . . .	297
18.2 Objectives . . . . .	297
18.3 Introduction . . . . .	297
18.4 Why does a project need explicit rules? . . . . .	298
18.5 How should I license my software? . . . . .	298
18.6 How should I license my data and reports? . . . . .	299
18.7 Why should I establish a code of conduct for my project? . . .	300
18.8 Why can I be a good ally for members of marginalized groups? .	301
18.9 How can I be a good ally for members of marginalized groups? .	302
18.10 Summary . . . . .	303
18.11 Exercises . . . . .	303
18.12 Key Points . . . . .	303
<b>19 Managing Backlog</b>	<b>305</b>

19.1 Questions . . . . .	305
19.2 Objectives . . . . .	305
19.3 Introduction . . . . .	305
19.4 How can I manage the work I still have to do? . . . . .	305
19.5 How can I write a good bug report? . . . . .	306
19.6 How can I use labels to organize work? . . . . .	307
19.7 How can I use labels to prioritize work? . . . . .	308
19.8 How can I use issues to enforce a workflow for a project? . . . . .	309
19.9 Summary . . . . .	310
19.10 Exercises . . . . .	310
19.11 Key Points . . . . .	310
<b>20 Teamwork</b>	<b>311</b>
20.1 Questions . . . . .	311
20.2 Objectives . . . . .	311
20.3 Introduction . . . . .	311
20.4 How can we run meetings more efficiently? . . . . .	311
20.5 What should we do when the meeting is over? . . . . .	313
20.6 How can we keep people from talking too much or too little? . . . . .	314
20.7 How should we run online meetings? . . . . .	314
20.8 What sort of people make teamwork hard? . . . . .	315
20.9 How should we handle conflict within the team? . . . . .	316
20.10 Summary . . . . .	318
20.11 Key Points . . . . .	318
<b>21 R Packaging</b>	<b>319</b>
21.1 Questions . . . . .	319
21.2 Objectives . . . . .	319
21.3 What's in an R package? . . . . .	319
21.4 What <i>is</i> a package, exactly? . . . . .	320
21.5 How do I create a package? . . . . .	321
21.6 How do I use a package while I'm creating it? . . . . .	324
21.7 How do I document a R package? . . . . .	324
21.8 What should I document? . . . . .	327
21.9 How do I manage package dependencies? . . . . .	328
21.10 How can I share my package via GitHub? . . . . .	330
21.11 How can I add data to a package? . . . . .	330
21.12 Summary . . . . .	333
21.13 Key Points . . . . .	333
<b>22 Python Packaging</b>	<b>335</b>
22.1 Questions . . . . .	335
22.2 Objectives . . . . .	335
22.3 Introduction . . . . .	335
22.4 How can I turn a set of Python source files into a module? . . . . .	336
22.5 How can I control what is executed during import and what isn't? . . . . .	337



22.6	How can I install a Python package? . . . . .	338
22.7	How can I create an installable Python package? . . . . .	338
22.8	How can I manage the source code of large packages? . . . . .	339
22.9	How can I distribute software packages that I have created? . . .	340
22.10	How can I manage the packages my projects need? . . . . .	342
22.11	How can I test package installation? . . . . .	344
22.12	Announcing Work . . . . .	344
22.13	Summary . . . . .	345
22.14	Exercises . . . . .	345
22.15	Key Points . . . . .	345
<b>23</b>	<b>Continuous Integration</b>	<b>347</b>
23.1	Questions . . . . .	347
23.2	Objectives . . . . .	347
23.3	Introduction . . . . .	347
23.4	The basics of setting up CI on a repository? . . . . .	348
23.5	Using Travis to test your software . . . . .	351
23.6	How can I use CI for non-testing purposes? . . . . .	352
23.7	Summary . . . . .	353
23.8	Exercises . . . . .	353
23.9	Key Points . . . . .	353
<b>24</b>	<b>Publishing</b>	<b>355</b>
24.1	Questions . . . . .	355
24.2	Objectives . . . . .	355
24.3	Introduction . . . . .	355
24.4	Identification . . . . .	356
24.5	Publishing a report . . . . .	356
24.6	Publishing data . . . . .	357
24.7	Publishing analysis software . . . . .	361
24.8	Publishing analysis scripts . . . . .	361
24.9	Summary . . . . .	364
24.10	Exercises . . . . .	364
24.11	Key Points . . . . .	365
<b>25</b>	<b>Finale</b>	<b>367</b>
25.1	Questions . . . . .	367
25.2	Objectives . . . . .	367
25.3	Introduction . . . . .	367
25.4	Programs are data. . . . .	367
25.5	Computers don't understand anything. . . . .	368
25.6	Programming is about creating and combining abstractions. . . .	368
25.7	Every redundancy in software is an abstraction trying to be born.	369
25.8	Create models for computers and views for human beings. . . . .	369
25.9	Paranoia makes us productive. . . . .	370

25.10 Things that don't change are easier to understand than things that do. . . . .	370
25.11 Better algorithms are better than better hardware. . . . .	371
25.12 Distributed is different. . . . .	371
25.13 Privacy, security, fairness, and responsibility can't be added after the fact. . . . .	371
25.14 Summary . . . . .	372
25.15 Exercises . . . . .	372
25.16 Key Points . . . . .	372
<b>A License</b>	<b>373</b>
<b>B Code of Conduct</b>	<b>375</b>
B.1 Our Standards . . . . .	375
B.2 Our Responsibilities . . . . .	376
B.3 Scope . . . . .	376
B.4 Enforcement . . . . .	376
B.5 Attribution . . . . .	376
<b>C Contributing</b>	<b>377</b>
<b>D Glossary</b>	<b>379</b>
<b>E Learning Objectives</b>	<b>385</b>
E.1 Novice R . . . . .	385
E.2 Novice Python . . . . .	385
E.3 RSE . . . . .	387
<b>F Key Points</b>	<b>389</b>
F.1 Novice R . . . . .	389
F.2 Novice Python . . . . .	389
F.3 RSE . . . . .	391
<b>G The Rules</b>	<b>393</b>
	<b>395</b>

# List of Tables

5.2 Table caption here. . . . . 77



# List of Figures

5.1	“Knit” button. . . . .	65
5.2	Steps to being more reproducible. Source DOI: 10.1038/d41586-018-05990-5 . . . . .	68
5.3	Caption. . . . .	69
5.4	Add your figure title here. . . . .	73
5.5	Add your figure title here. . . . .	75
5.6	Add your figure title here. . . . .	75
9.1	Rich rendering of the docstring in Spyder . . . . .	141
13.1	Configuration Concept Map . . . . .	243
14.1	Automation Concept Map . . . . .	260
16.1	Number Spacing . . . . .	279
17.1	Project Layout . . . . .	291
19.1	Issue State Transitions . . . . .	309
21.1	Package Distribution . . . . .	329
21.2	R Packaging Concept Map . . . . .	334



# Chapter 1

## Overview

It’s still magic even if you know how it’s done.

– Terry Pratchett

FIXME: general introduction

- Be able to do the steps in Yenni et al. (2019).

### 1.1 Why isn’t all of this normal already?

Nobody argues that research should be irreproducible or unsustainable, but “not against it” and actively supporting it are very different things. Academia doesn’t yet know how to reward people for writing useful software, so while you may be thanked, the extra effort you put in may not translate into job security or decent pay.

And some people still argue against openness, Being open is a big step toward a (non-academic) career path, which is where approximately 80% of PhDs go, and for those staying in academia, open work is cited more often than closed (FIXME: citation). However, some people still worry that if they make their data and code generally available, someone else will use it and publish a result they have come up with themselves. This is almost unheard of in practice, but that doesn’t stop people using it as a boogeyman.

Other people are afraid of looking foolish or incompetent by sharing code that might contain bugs. This isn’t just impostor syndrome: members of marginalized groups are frequently judged more harshly than others (FIXME: CITE).

## 1.2 Acknowledgments

This book owes its existence to the hundreds of researchers we met through Software Carpentry and Data Carpentry. We are also grateful to Insight Data Science for sponsoring the early stages of this work, to everyone who has contributed, and to:

- *Practical Computing for Biologists* Haddock and Dunn (2010)
- “A Quick Guide to Organizing Computational Biology Projects” Noble (2009)
- “Ten Simple Rules for Making Research Software More Robust” Taschuk and Wilson (2017)
- “Best Practices for Scientific Computing” Wilson et al. (2014)
- “Good Enough Practices in Scientific Computing” Wilson et al. (2017)



## Chapter 2

# Novice Goals

This outline describes the questions that the novice courses on R and Python will answer. The advanced course can then assume that learners have hands-on experience with these topics but nothing more.

## 2.1 Personas

### 2.1.1 Anya

**Anya** is a professor of neuropsychology who is responsible for teaching her department's introduction to statistics to 1100 first-year students every year. (Students complain that the Stats department's introductory course is too theoretical and requires more programming knowledge than they have.) When she finds time for it, her research focuses on color perception in infants.

Over the past nine years, Anya has designed and run a dozen experiments on 50-100 infant subjects each and analyzed the results using SPSS and more recently R (which she taught herself during a sabbatical). She has never taken a programming course, and suffers from impostor syndrome when talking to colleagues who are using things like GitHub and R Markdown.

Anya would like to figure out how to use R to teach her intro stats course, which currently uses a mixture of Excel and SPSS. She would like to learn more about time series analysis to support her research, and about tools like Git and R Markdown.

This guide has modular lessons and exercises that she can adapt to use in her course, and suggestions for how to make learning interactive with a large class size. She also finds helpful instructions for applying time series analysis to data using R.

### 2.1.2 Exton

**Exton** taught business at a community college before joining a friend’s startup, and now does community management for a company that builds healthcare software. He still teaches Marketing 101 every year to help people with backgrounds like his.

Exton uses Excel to keep track of who is registered for webinars, workshops, and training sessions. Some of these spreadsheets are created from CSV files produced by a web-scraping script a summer intern wrote for him a couple of years ago. Exton doesn’t think of himself as a programmer, but spends hours creating complicated lookup tables in multi-sheet spreadsheets to help him figure out how many webinar attendees turn into community contributors, who answers forum posts most frequently, and so on.

Exton knows there are better ways to do what he’s doing, but feels overwhelmed by the flood of blog posts, tweets, and “helpful” recommendations he receives from members of the company’s engineering team. He wants someone to tell him where he should start and how long it will take whatever he learns to pay off.

Exton finds ‘Merely Useful’ after some Googling, and sees an example of data analysis with spreadsheet data that looks really similar to what he’s trying to do. He carefully works through that particular example, then goes back and works through some of the earlier material in the book. He can tell that it won’t take long to get this to work with his data.

### 2.1.3 Irwin

**Irwin**, 18, is five months into an undergraduate degree in urban planning. He’s read lots of gushing articles in *Wired* about data science, and was excited by the prospect of learning how to do it, but dropped his CS 101 course after six weeks because nothing made sense. (His university’s computer science department uses Haskell as an introductory programming language...) He is doing better in Anya’s course (which he is taking as an elective) but still spends most of his time copying, pasting, and swearing.

Irwin did well in his high school math classes, and built himself a home page with HTML and CSS in a weekend workshop in grade 11. He knows how to do simple calculations in Excel, has accounts on nine different social media sites, and attends all of his morning classes online.

Anya mentions this guide in one of her classes, and Irwin downloads the PDF to read on the bus. He loves the examples that use urban data, and right away he has tons of ideas about where to get more cool data to analyze. His urban data science blog is already taking shape in his head.

### 2.1.4 Camilla

**Camilla** recently started a job as an assistant professor. Her department (Medieval Studies) is trying to develop a digital humanities data-science-heavy undergraduate program, and the undergraduate chair thinks that Camilla has the most programming experience in the department and has asked her to develop an introduction to programming course for humanities students.

Camilla has dabbled in natural language processing and has learned Python over the course of her previous work, but she has no experience teaching programming and she's not sure what the best way is to teach beginners. She doesn't want to start from scratch to create a course out of nothing. She also isn't sure which programming language the new program should focus on.

She finds 'Merely Useful' and feels relieved: she can pretty much use the book as-is for her course. She looks up the examples of text and image analysis and compares how both R and Python approach those kinds of data to help her make a decision about which language to teach.

### 2.1.5 Jordan

**Jordan** is a third-year undergraduate student in ecology. Two months ago she started working part-time for a professor in her department, and she's beginning to collect and analyze data from her own experiments with fruit flies. Her professor has asked her to learn R to do her analysis and suggested that she sign up for the introduction to quantitative data analysis in R course that the ecology department offers. The course is just starting, and it uses 'Merely Useful' as the textbook.

Jordan can't wait to apply her new programming knowledge to her data, so she starts reading ahead and trying to use her own data in some of the book's examples. As she works through examples, she realizes that she'll need to change a few things about how she records her data in spreadsheets so that it will be easier to analyze in R.

## 2.2 Getting started

- What are the different ways I can interact with software?
  - console
  - scripts
- How can I find and view help?
  - In the IDE
  - Stack Overflow
- How can I inspect data while I'm working on it?
  - table viewers

## 2.3 Data manipulation

- How can I read tabular data into a program?
  - what CSV is, where it comes from, and why people use it
  - reading files
- How can I select subsets of my data?
  - select
  - filter
  - arrange
  - Boolean conditions
- How can I calculate new values?
  - mutate
  - ifelse
- How can I tell what’s gone wrong in my programs?
  - reading error messages
  - the difference between syntax and runtime errors
- How can I operate on subsets of my data?
  - group
  - summarize
  - split-apply-combine
- How can I work with two or more datasets?
  - join
- How can I save my results?
  - writing files
- What *isn’t* included?
  - anything other than reasonably tidy tabular data
  - map
  - loops and conditionals

## 2.4 Plotting

- Why plot?
  - summary statistics can mislead
  - Anscombe’s Quartet and the DataSaurus dozen
- What are the core elements of every plot?
  - data
  - geometric objects
  - aesthetic mapping
- How can I create different kinds of plots?
  - scatter plot
  - line plot
  - histogram
  - bar plot
  - which to use when

- How can I plot multiple datasets at once?
  - grouping
  - faceting
- How can I make misleading plots?
  - showing a single central tendency data point instead of the individual observations
  - saturated plots instead of for example violins or 2D histograms
  - picking unreasonable axes limits to intentionally misrepresent the underlying data
  - not using perceptually uniform colormaps to indicate quantities
  - not thinking about color blindness
- What *isn't* included?
  - outliers
  - interactive plots
  - maps
  - 3D visualization

## 2.5 Development

- How can I make my own functions?
  - declaring functions
  - declaring parameters
  - default values
  - common conventions
- How can I make my programs tell me that something has gone wrong?
  - validation (did we build the right thing) vs. verification (did we build the thing correctly)
  - assertions for sanity checks
- How can I ask for help?
  - creating a reproducible example (reprex)
- How do I install software?
  - what *is* a package?
  - package manager
- What *isn't* included?
  - code browsers, multiple cursors, and other fancy IDE tricks
  - virtual environments
  - debuggers

## 2.6 Data analysis

- How can I represent and manage missing values?
  - NA

- How can I get a feel for my data?
  - summary statistics
- How can I create a simple model of my data?
  - formulas
  - linear regression
    - \* adding a best fit straight line on a scatterplot
    - \* understanding what the error bands on a “best fit” straight line mean
  - k-means cluster analysis
  - frame these as exploratory tools for revealing structure in the data, rather than modelling or inferential tools
- How can I put people at risk?
  - algorithmic bias
  - de-anonymization
- What *isn't* included?
  - statistical tests
  - multiple linear regression
  - anything with “machine learning” in its name

## 2.7 Version control

- What is a version control system?
  - a smarter kind of backup
- What goes where and why?
  - local vs. remote storage (physically)
  - local vs. remote storage (ethical/privacy issues)
- How do I track my work locally?
  - diff
  - add
  - commit
  - log
- How do I view or recover an old version of a file?
  - diff
  - checkout
- How do I save work remotely?
  - push and pull
- How do I manage conflicts?
  - merge
- What *isn't* included?
  - forking
  - branching
  - pull requests
  - git reflow –substantive –single-afferent-cycle –ia-ia-rebase-fhtagn ...

## 2.8 Publishing

- How do static websites work?
  - URLs
  - servers
  - request/response cycle
  - pages
- How do I create a simple HTML page?
  - head/body
  - basic elements
  - images
  - links
  - relative vs. absolute paths
- How can I create a simple website?
  - GitHub pages
- How can I give pages a standard appearance?
  - layouts
- How can I avoid writing all those tags?
  - Markdown
- How can I share values between pages?
  - flat per-site and per-page configuration
  - variable expansion
- What *isn't* included?
  - templating
  - filters
  - inclusions

## 2.9 Reproducibility

- How can I make programs easy to read?
  - coding style
  - linters
  - documentation
- How can I make programs easy to re-use?
  - Taschuk's Rules
- How can I combine explanations, code, and results?
  - notebooks
- Where does stuff actually live on my computer?
  - directory structure on Windows and Unix
  - absolute vs. relative paths
  - significance of the working directory
  - data on disk vs. data in memory
- How should I organize my projects?
  - Noble's Rules

- RStudio projects
- How should I keep track of my data?
  - simple manifests
- What *isn't* included?
  - build tools (Make and its kin)
  - continuous integration
  - documentation generators

## 2.10 Collaboration

- What kinds of licenses are there?
  - open vs. closed
  - copyright
- Who gets to decide what license to use?
  - it depends...
- What license should I use for my publications?
  - CC-something
- What license should I use for my software?
  - MIT/BSD vs. GPL
- What license should I use for my data?
  - CC0
- How should I identify myself and my work?
  - DOIs
  - ORCIDs
- How do I credit someone else's code?
  - citing packages, citing something from GitHub, giving credit for someone's answer on StackOverflow...
- What's the difference between open and welcoming?
  - evidence for systematic exclusion
  - mechanics of exclusion
- How can I help create a level playing field?
  - what's wrong with deficit models
  - allyship, advocacy, and sponsorship
  - Code of Conduct (remove negatives)
  - curb cuts (adding positives for some people helps everyone else too)
- What *isn't* included?
  - how to run a meeting
  - community management
  - mental health
  - assessment of this course



**Part I**

**Novice R Material**



## Chapter 3

# Introduction

**Note:** we will include this introduction in the novice Python material as well.

FIXME: general introduction.

### 3.1 Who are these lessons for?

#### 3.1.1 Exton Excel

1. Exton taught business at a community college for several years, and now does community management for an event management company. He still teaches Marketing 101 every year to help people with backgrounds like his.
2. Exton uses Excel to keep track of who is registered for webinars, workshops, and training sessions. He doesn't think of himself as a programmer, but spends hours creating complicated lookup tables to figure out how many webinar attendees turn into community contributors, who answers forum posts most frequently, and so on.
3. Exton knows there are better ways to do what he's doing, but feels overwhelmed by the blog posts, tweets, and "helpful" recommendations from the company's engineering team.
4. Exton is a single parent; the one evening a week he spends teaching is the only out-of-work time he can take away from family responsibilities.

### 3.1.2 Nina Newbie

1. Nina is 18 years old and in the first year of an undergraduate degree in urban planning. She has read lots of gushing articles about data science, and was excited by the prospect of learning how to do it, but dropped her CS 101 course after six weeks because nothing made sense. She is doing better in her intro to statistics (which uses a little bit of R), but still spends most of her time copying, pasting, and swearing.
2. Nina did well in her high school math classes, and built herself a home page with HTML and CSS in a weekend workshop in grade 11. She has accounts on nine different social media site, and attends all of her morning classes online.
3. Nina wants self-paced tutorials with practice exercises, plus forums where she can ask for help.
4. After a few bruising conversations with CS majors, Nina is reluctant to reveal how little she knows about programming: she would rather get a low grade and blame it on partying than let her classmates see that she is floundering.

## 3.2 What will these lessons teach you?

During this course, you will learn how to:

- Write programs in R that read data, clean it up, perform simple statistical analyses on it.
- Build visualizations to help you understand your data and communicate your findings.
- Find and install R packages to help you do these things.
- Create software that other people can understand and re-run.
- Track your work in version control using Git and GitHub.
- Publish your results on a web site using R Markdown and GitHub Pages.
- Make your data, software, and reports citable using ORCIDs and DOIs, and cite the work of others.
- Select open source and Creative Commons licenses that allow you and others to share data, software, and reports.
- Be an active participant in open, inclusive projects.

### 3.2.1 What do you need to have and know to start?

This course assumes that you:

- Have a laptop that you can install software on, or access to a web-based programming system like `rstudio.cloud`.
- Know what mean and variance are.

- Are willing to invest about 30 hours in reading or lectures and another 100 hours doing practice exercises.

We have tried to make these lessons accessible to people with visual or motor challenges, but recognize that some parts (particularly data visualization) may still be difficult. We welcome suggestions for improvements.

### 3.3 What examples will we use?

FIXME: introduce running examples

### 3.4 What's the big picture? (`#novice-r-intro-bigpicture`)

We now swim in a sea of data and generate more each day. That data can help us understand the world, but it can also be used to manipulate us and invade our privacy. Learning how to analyze data will help you do the former and guard against the latter.

This course is therefore about *people*, *programs*, and *data*. Data can live in three places (FIXME: diagram)

1. In the computer's memory. It has to be here for the computer to use it, but when a program stops running or the computer is shut down, the contents of memory evaporate.
2. On the computer's hard drive. This is much larger than memory—terabytes instead of gigabytes—and its contents are organized into *files* and *directories* (also called *folders*). What's on the hard drive stays there even when programs aren't running or the computer is switched off. A program must *read* data from files into memory to work with it, and *write* data to files to save it permanently.
3. On some other computer on the network. “The cloud” and “the web” are just other people's computers; if we want to use data that's on the other side of a URL, we need to download it (i.e., copy it to our hard drive) or read it directly into memory (which is called streaming).

Programs also live in three places (FIXME: enhanced diagram)

1. In the computer's memory. A program has to be in memory for the computer to run it.
2. On the computer's hard drive. A program is a file like any other. Instead of containing text, pixels, or CO2 measurements, it contains instructions. In order for the computer to run it, those instructions have to be copied into

memory. (This is part of what your computer is doing when it launches an application.) Your program typically uses other pieces of software called packages or libraries that provide common operations like searching in text, changing the colors of pixels, or calculating averages. When your program is loaded into memory, your computer also loads those packages.

3. On some other computer on the network. R and Python both have catalogs of packages that people have written and shared. In order to use one of these, you must install it on your computer by copying its files from the catalog site to your hard drive. There's usually more to this than simply copying one file, so R and Python both come with tools to help you find, install, and manage packages.

Finally, we come to people (FIXME: enhanced diagram)

1. This course starts by teaching you how to write programs that run on your computer, analyze data that is on your computer or on the web, and use packages written by other people.
2. It will also teach you how to write programs that your collaborators can understand and run. (One of those collaborators is your future self, who will be grateful three months from now that you did the right thing today.) These skills will make you a productive member of a small team, and this course will also explain how to make such teams open and inclusive.
3. The third group of people includes collaborators or reviewers who want to reproduce your work. This course will show you how to publish your results on the web and how to get credit for your work and give it to others.

## Chapter 4

# Practice

We have covered a lot in the last few lessons, so this one presents some practice exercises to ground what we have learned and introduce a few more commonly-used functions.

### 4.1 Working with a single tidy table

1. Load the tidyverse collection of package and the `here` package for constructing paths:

```
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.2.1 --
## v ggplot2 3.2.1      v purrr   0.3.2
## v tibble  2.1.3      v dplyr  0.8.3
## v tidyr   1.0.0      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.4.0

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()

library(here)
```

```
## here() starts at /home/travis/build/merely-useful/merely-useful.github.io
```

2. Use `here::here` to construct a path to a file and `readr::read_csv` to read that file:

```
path <- here::here("data", "person.csv")
person <- readr::read_csv(path)
```

```
## Parsed with column specification:
## cols(
##   person_id = col_character(),
##   personal_name = col_character(),
##   family_name = col_character()
## )
```

```
person
```

```
## # A tibble: 5 x 3
##   person_id personal_name family_name
##   <chr>      <chr>      <chr>
## 1 dyer      William      Dyer
## 2 pb       Frank       Pabodie
## 3 lake     Anderson    Lake
## 4 roe      Valentina   Roerich
## 5 danforth Frank       Danforth
```

Read *survey/site.csv*.

3. Count rows and columns using `nrow` and `ncol`:

```
nrow(person)
```

```
## [1] 5
```

```
ncol(person)
```

```
## [1] 3
```

How many rows and columns are in the site data?

4. Format strings using `glue::glue`:

```
print(glue::glue("person has {nrow(person)} rows and {ncol(person)} columns"))
```

```
## person has 5 rows and 3 columns
```

Print a nicely-formatted summary of the number of rows and columns in the site data.

5. Use `colnames` to get the names of columns and `paste` to join strings together:

```
print(glue::glue("person columns are {paste(colnames(person), collapse = ' ')}"))
```

```
## person columns are person_id personal_name family_name
```

Print a nicely-formatted summary of the names of the columns in the site data.

6. Use `dplyr::select` to create a new table with a subset of columns by name:



```
dplyr::select(person, family_name, personal_name)
```

```
## # A tibble: 5 x 2
##   family_name personal_name
##   <chr>        <chr>
## 1 Dyer         William
## 2 Pabodie      Frank
## 3 Lake         Anderson
## 4 Roerich      Valentina
## 5 Danforth     Frank
```

Create a table with just the latitudes and longitudes of sites.

7. Use `dplyr::filter` to create a new table with a subset of rows by values:

```
dplyr::filter(person, family_name < "M")
```

```
## # A tibble: 3 x 3
##   person_id personal_name family_name
##   <chr>      <chr>        <chr>
## 1 dyer      William      Dyer
## 2 lake      Anderson    Lake
## 3 danforth  Frank       Danforth
```

Create a table with only sites south of -48 degrees.

8. Use the pipe operator `%>%` to combine operations:

```
person %>%
  dplyr::select(family_name, personal_name) %>%
  dplyr::filter(family_name < "M")
```

```
## # A tibble: 3 x 2
##   family_name personal_name
##   <chr>        <chr>
## 1 Dyer         William
## 2 Lake         Anderson
## 3 Danforth     Frank
```

Create a table with only the latitudes and longitudes of sites south of -48 degrees.

9. Use `dplyr::mutate` to create a new column with calculated values and `stringr::str_length` to calculate string length:

```
person %>%
  dplyr::mutate(name_length = stringr::str_length(family_name))
```

```
## # A tibble: 5 x 4
##   person_id personal_name family_name name_length
##   <chr>      <chr>        <chr>         <int>
```

```
## 1 dyer      William      Dyer      4
## 2 pb       Frank       Pabodie    7
## 3 lake     Anderson    Lake       4
## 4 roe      Valentina   Roerich    7
## 5 danforth Frank       Danforth    8
```

Look at the help for the built-in function `round` and then use it to create a table with latitudes and longitudes rounded to integers.

10. Use `dplyr::arrange` to order rows and (optionally) `dplyr::desc` to impose descending order:

```
person %>%
  dplyr::mutate(name_length = stringr::str_length(family_name)) %>%
  dplyr::arrange(dplyr::desc(name_length))

## # A tibble: 5 x 4
##   person_id personal_name family_name name_length
##   <chr>      <chr>        <chr>         <int>
## 1 danforth  Frank          Danforth         8
## 2 pb       Frank          Pabodie         7
## 3 roe      Valentina     Roerich         7
## 4 dyer     William       Dyer            4
## 5 lake     Anderson     Lake            4
```

Create a table sorted by decreasing longitude (i.e., most negative longitude last).

## 4.2 Working with grouped data

1. Read `survey/measurements.csv` and look at the data with `View`:

```
measurements <- readr::read_csv(here::here("data", "measurements.csv"))

## Parsed with column specification:
## cols(
##   visit_id = col_double(),
##   visitor = col_character(),
##   quantity = col_character(),
##   reading = col_double()
## )
View(measurements)
```

2. Find rows where `reading` is not NA, save as `cleaned`, and report how many rows were removed:

```
cleaned <- measurements %>%
  dplyr::filter(!is.na(reading))
```

```
nrow(measurements) - nrow(cleaned)
```

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

*Rewrite the filter expression to select rows where the visitor and quantity are not NA either and report the total number of rows removed.*

3. Group measurements by quantity measured and count the number of each (the column is named n automatically):

```
cleaned %>%
  dplyr::group_by(quantity) %>%
  dplyr::count()
```

```
## # A tibble: 3 x 2
## # Groups:   quantity [3]
##   quantity     n
##   <chr>    <int>
## 1 rad         8
## 2 sal         7
## 3 temp        3
```

*Group by person and quantity measured.*

4. Find the minimum, average, and maximum for each quantity:

```
cleaned %>%
  dplyr::group_by(quantity) %>%
  dplyr::summarize(low = min(reading), mid = mean(reading), high = max(reading))
```

```
## # A tibble: 3 x 4
##   quantity    low    mid  high
##   <chr>    <dbl> <dbl> <dbl>
## 1 rad      1.46  6.56  11.2
## 2 sal      0.05  9.24  41.6
## 3 temp   -21.5 -18.7  -16
```

*Look at the range for each combination of person and quantity.*

5. Rescale salinity measurements that are greater than 1:

```
cleaned <- cleaned %>%
  dplyr::mutate(reading = ifelse(quantity == 'sal' & reading > 1.0, reading/100, reading))
cleaned
```

```
## # A tibble: 18 x 4
##   visit_id visitor quantity reading
##     <dbl> <chr>    <chr>    <dbl>
## 1     619 dyer    rad      9.82
## 2     619 dyer    sal       0.13
## 3     622 dyer    rad       7.8
```

```
## 4      622 dyer    sal      0.09
## 5      734 pb     rad      8.41
## 6      734 lake   sal      0.05
## 7      734 pb     temp    -21.5
## 8      735 pb     rad      7.22
## 9      751 pb     rad      4.35
## 10     751 pb     temp    -18.5
## 11     752 lake   rad      2.19
## 12     752 lake   sal      0.09
## 13     752 lake   temp    -16
## 14     752 roe    sal      0.416
## 15     837 lake   rad      1.46
## 16     837 lake   sal      0.21
## 17     837 roe    sal      0.225
## 18     844 roe    rad      11.2
```

Do the same calculation use *case\_when*.

6. Read `visited.csv`, drop the NAs and store in `visits`. Use `anti_join()` to find the measurements in `cleaned` that don't have matches in `visits`:

```
visits <- readr::read_csv(here::here("data", "visited.csv")) %>%
  dplyr::filter(!is.na(visit_date))
```

```
## Parsed with column specification:
## cols(
##   visit_id = col_double(),
##   site_id = col_character(),
##   visit_date = col_date(format = "")
## )
```

```
cleaned %>% anti_join(visits)
```

```
## Joining, by = "visit_id"
```

```
## # A tibble: 4 x 4
##   visit_id visitor quantity reading
##   <dbl> <chr>    <chr>    <dbl>
## 1     752 lake    rad      2.19
## 2     752 lake    sal      0.09
## 3     752 lake    temp    -16
## 4     752 roe    sal      0.416
```

Are there any sites in *visits* that don't have matches in *cleaned*?

7. Join `visits` with the cleaned-up table of readings:

```
cleaned <- visits %>%
  dplyr::inner_join(cleaned, by = c("visit_id" = "visit_id"))
cleaned
```

```
## # A tibble: 14 x 6
##   visit_id site_id visit_date visitor quantity reading
##   <dbl> <chr> <date> <chr> <chr> <dbl>
## 1     619 DR-1  1927-02-08 dyer    rad     9.82
## 2     619 DR-1  1927-02-08 dyer    sal     0.13
## 3     622 DR-1  1927-02-10 dyer    rad     7.8
## 4     622 DR-1  1927-02-10 dyer    sal     0.09
## 5     734 DR-3  1930-01-07 pb     rad     8.41
## 6     734 DR-3  1930-01-07 lake    sal     0.05
## 7     734 DR-3  1930-01-07 pb     temp    -21.5
## 8     735 DR-3  1930-01-12 pb     rad     7.22
## 9     751 DR-3  1930-02-26 pb     rad     4.35
## 10    751 DR-3  1930-02-26 pb     temp    -18.5
## 11    837 MSK-4  1932-01-14 lake    rad     1.46
## 12    837 MSK-4  1932-01-14 lake    sal     0.21
## 13    837 MSK-4  1932-01-14 roe    sal     0.225
## 14    844 DR-1  1932-03-22 roe    rad     11.2
```

Join *visited.csv* with *site.csv* to get (date, latitude, longitude) triples for site visits.

7. Find the dates of the highest radiation reading at each site:

```
cleaned %>%
  dplyr::filter(quantity == "rad") %>%
  dplyr::group_by(site_id) %>%
  dplyr::mutate(max_rad = max(reading)) %>%
  dplyr::filter(reading == max_rad)

## # A tibble: 3 x 7
## # Groups:   site_id [3]
##   visit_id site_id visit_date visitor quantity reading max_rad
##   <dbl> <chr> <date> <chr> <chr> <dbl> <dbl>
## 1     734 DR-3  1930-01-07 pb     rad     8.41  8.41
## 2     837 MSK-4  1932-01-14 lake    rad     1.46  1.46
## 3     844 DR-1  1932-03-22 roe    rad    11.2  11.2
```

Another way to do it:

```
cleaned %>%
  dplyr::filter(quantity == "rad") %>%
  dplyr::group_by(site_id) %>%
  dplyr::top_n(1, reading) %>%
  dplyr::select(site_id, visit_date, reading)

## # A tibble: 3 x 3
## # Groups:   site_id [3]
##   site_id visit_date reading
```

```
##   <chr>   <date>       <dbl>
## 1 DR-3    1930-01-07    8.41
## 2 MSK-4    1932-01-14    1.46
## 3 DR-1    1932-03-22   11.2
```

Explain why this *doesn't* work.

```
cleaned %>%
  dplyr::filter(quantity == "rad") %>%
  dplyr::group_by(site_id) %>%
  dplyr::summarize(max_rad = max(reading)) %>%
  dplyr::ungroup() %>%
  dplyr::filter(reading == max_rad)
```

```
## Error: object 'reading' not found
```

8. Normalize radiation against the highest radiation seen per site:

```
cleaned %>%
  dplyr::filter(quantity == "rad") %>%
  dplyr::group_by(site_id) %>%
  dplyr::mutate(
    max_rad = max(reading),
    frac_rad = reading / max_rad) %>%
  dplyr::select(visit_id, site_id, visit_date, frac_rad)
```

```
## # A tibble: 7 x 4
## # Groups:   site_id [3]
##   visit_id site_id visit_date frac_rad
##       <dbl> <chr>   <date>       <dbl>
## 1      619 DR-1    1927-02-08    0.873
## 2      622 DR-1    1927-02-10    0.693
## 3      734 DR-3    1930-01-07     1
## 4      735 DR-3    1930-01-12    0.859
## 5      751 DR-3    1930-02-26    0.517
## 6      837 MSK-4    1932-01-14     1
## 7      844 DR-1    1932-03-22     1
```

Normalize salinity against mean salinity by site.

9. Find stepwise change in radiation per site by date:

```
cleaned %>%
  dplyr::filter(quantity == "rad") %>%
  dplyr::group_by(site_id) %>%
  dplyr::mutate(delta_rad = reading - dplyr::lag(reading)) %>%
  dplyr::arrange(site_id, visit_date)
```

```
## # A tibble: 7 x 7
## # Groups:   site_id [3]
```

```
##   visit_id site_id visit_date visitor quantity reading delta_rad
##   <dbl> <chr>   <date>      <chr>   <chr>      <dbl>   <dbl>
## 1      619 DR-1   1927-02-08 dyer    rad        9.82    NA
## 2      622 DR-1   1927-02-10 dyer    rad        7.8     -2.02
## 3      844 DR-1   1932-03-22 roe     rad       11.2     3.45
## 4      734 DR-3   1930-01-07 pb      rad        8.41    NA
## 5      735 DR-3   1930-01-12 pb      rad        7.22   -1.19
## 6      751 DR-3   1930-02-26 pb      rad        4.35   -2.87
## 7      837 MSK-4  1932-01-14 lake    rad        1.46    NA
```

*Find length of time between visits by site.*

10. Find sites that experience any stepwise increase in radiation between visits:

```
cleaned %>%
  dplyr::filter(quantity == "rad") %>%
  dplyr::group_by(site_id) %>%
  dplyr::mutate(delta_rad = reading - dplyr::lag(reading)) %>%
  dplyr::filter(!is.na(delta_rad)) %>%
  dplyr::summarize(any_increase = any(delta_rad > 0)) %>%
  dplyr::filter(any_increase)
```

```
## # A tibble: 1 x 2
##   site_id any_increase
##   <chr>   <lgl>
## 1 DR-1    TRUE
```

*Find sites with visits more than one year apart.*

## 4.3 Creating charts

We will use data on the mass and home range area (HRA) of various species from:

Tamburello N, Côté IM, Dulvy NK (2015) Data from: Energy and the scaling of animal space use. Dryad Digital Repository. <https://doi.org/10.5061/dryad.q5j65>

```
hra <- readr::read_csv(here::here("data", "home-range-database.csv"))
```

```
## Parsed with column specification:
## cols(
##   .default = col_character(),
##   mean.mass.g = col_double(),
##   log10.mass = col_double(),
##   mean.hra.m2 = col_double(),
```

```
## log10.hra = col_double(),
## preymass = col_double(),
## log10.preymass = col_double(),
## PPMR = col_double()
## )
```

```
## See spec(...) for full column specifications.
```

```
head(hra)
```

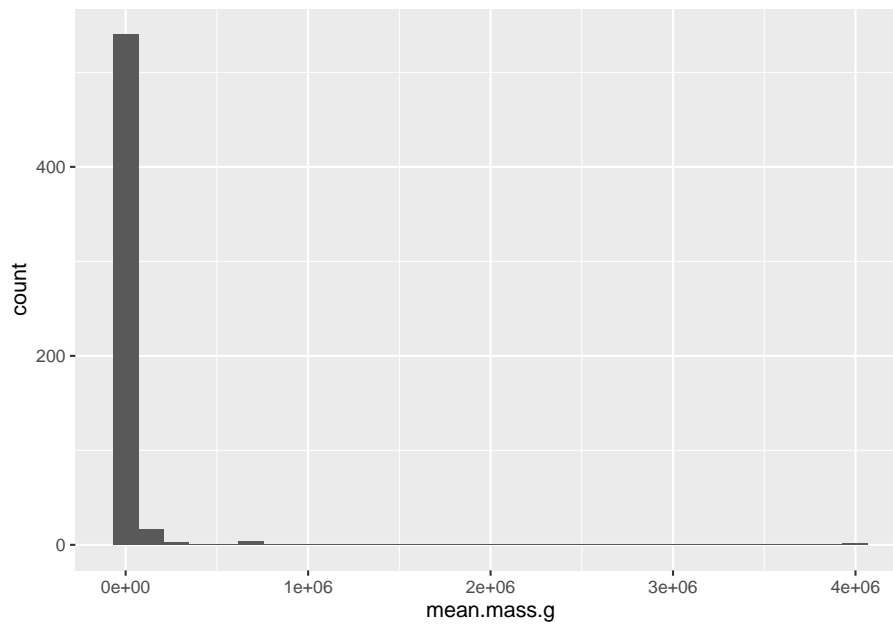
```
## # A tibble: 6 x 24
##   taxon common.name class order family genus species primarymethod N
##   <chr> <chr>         <chr> <chr> <chr> <chr> <chr> <chr> <chr>
## 1 lake~ american e~ acti~ angu~ angui~ angu~ rostra~ telemetry 16
## 2 rive~ blacktail ~ acti~ cypr~ catos~ moxo~ poecil~ mark-recaptu~ <NA>
## 3 rive~ central st~ acti~ cypr~ cypri~ camp~ anomal~ mark-recaptu~ 20
## 4 rive~ rosyside d~ acti~ cypr~ cypri~ clin~ fundul~ mark-recaptu~ 26
## 5 rive~ longnose d~ acti~ cypr~ cypri~ rhin~ catara~ mark-recaptu~ 17
## 6 rive~ muskellunge acti~ esoc~ esoci~ esox masqui~ telemetry 5
## # ... with 15 more variables: mean.mass.g <dbl>, log10.mass <dbl>,
## #   alternative.mass.reference <chr>, mean.hra.m2 <dbl>, log10.hra <dbl>,
## #   hra.reference <chr>, realm <chr>, thermoregulation <chr>,
## #   locomotion <chr>, trophic.guild <chr>, dimension <chr>,
## #   preymass <dbl>, log10.preymass <dbl>, PPMR <dbl>,
## #   prey.size.reference <chr>
```

1. Look at how mass is distributed:

```
ggplot2::ggplot(hra, mapping = aes(x = mean.mass.g)) +
  ggplot2::geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

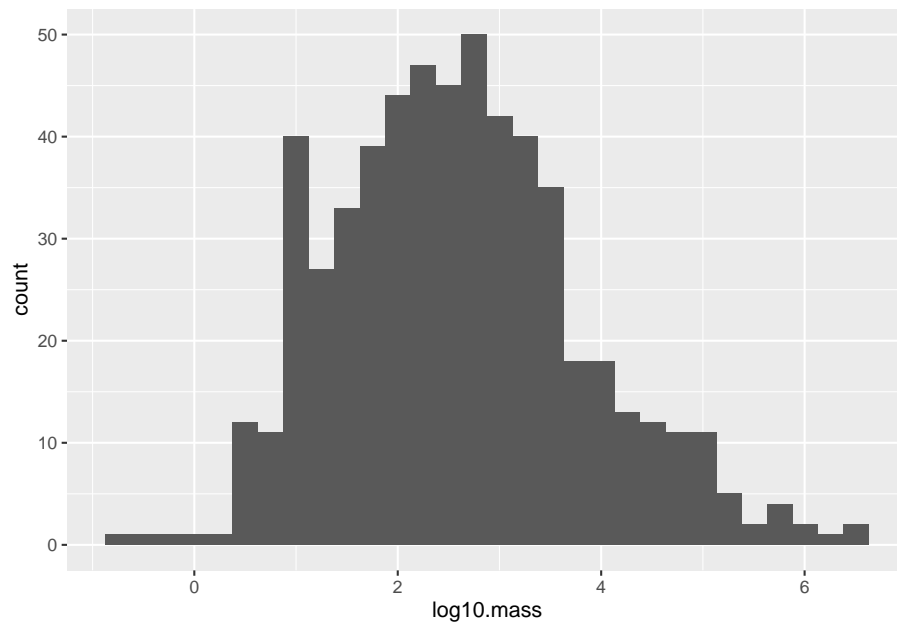




Try again with `log10.mass`:

```
ggplot2::ggplot(hra, mapping = aes(x = log10.mass)) +  
  ggplot2::geom_histogram()
```

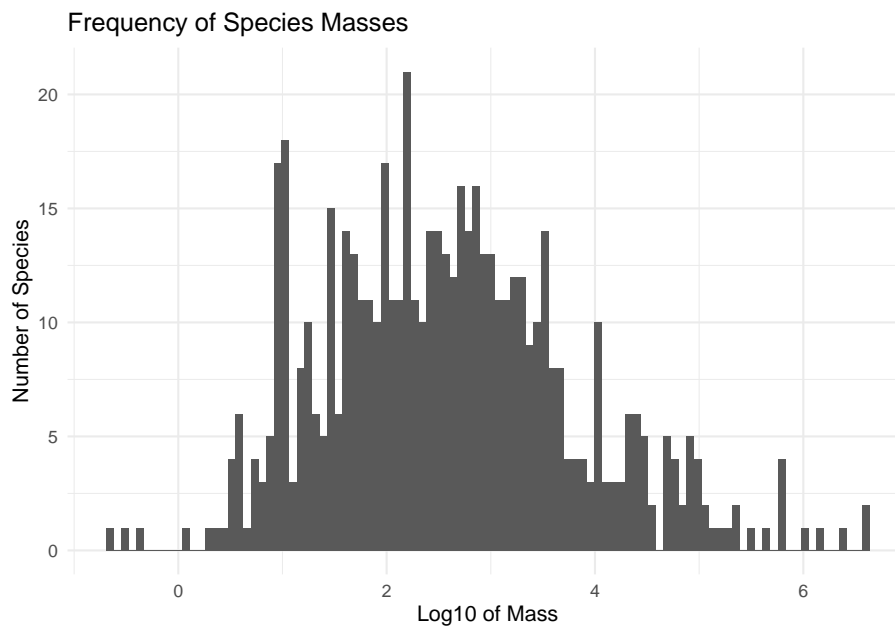
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Create histograms showing the distribution of home range area using linear and log scales.

2. Change the visual appearance of a chart:

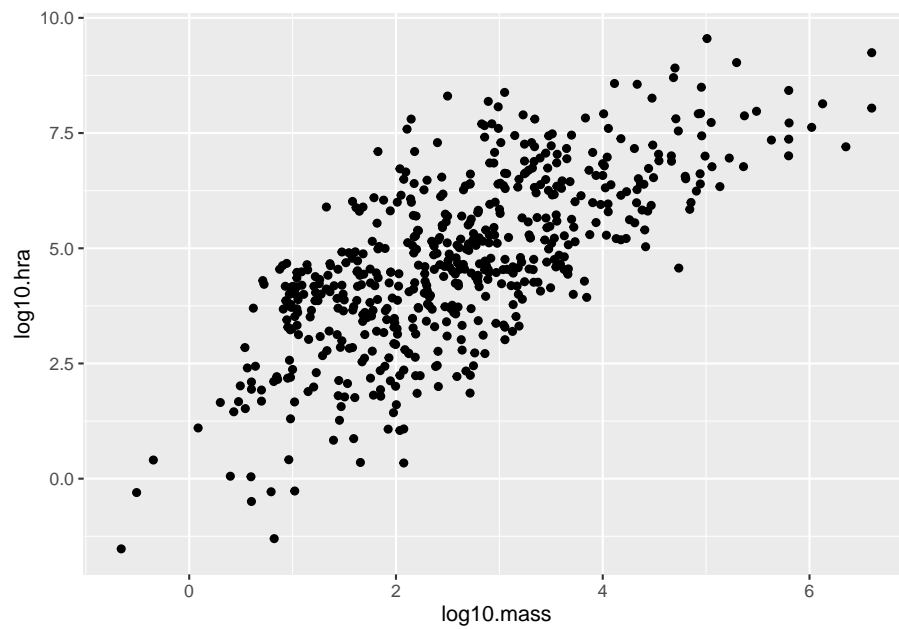
```
ggplot2::ggplot(hra, mapping = aes(x = log10.mass)) +  
  ggplot2::geom_histogram(bins = 100) +  
  ggplot2::ggtitle("Frequency of Species Masses") +  
  ggplot2::xlab("Log10 of Mass") +  
  ggplot2::ylab("Number of Species") +  
  ggplot2::theme_minimal()
```



*Show the distribution of home range areas with a dark background.*

3. Create a scatterplot showing the relationship between mass and home range area:

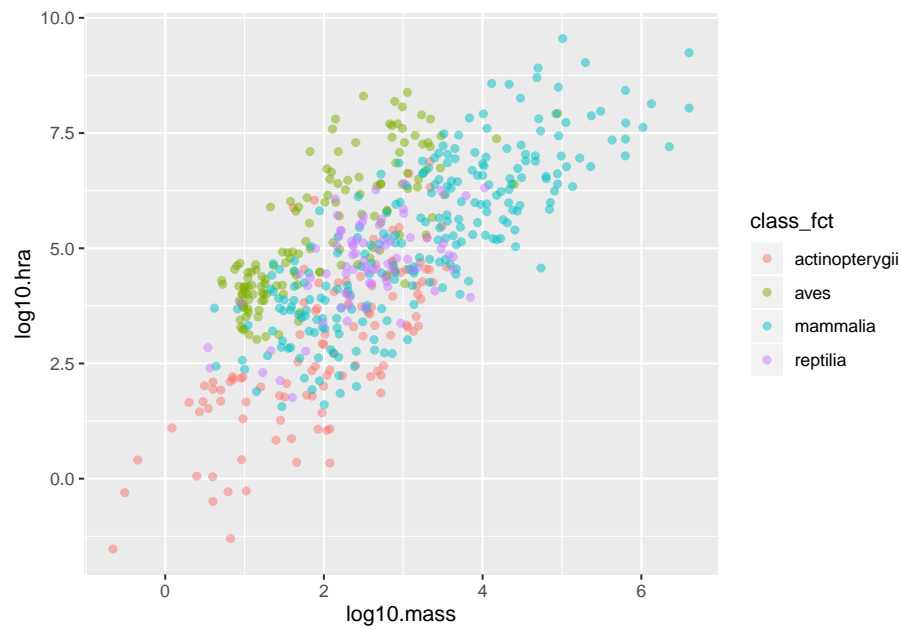
```
ggplot2::ggplot(hra, mapping = aes(x = log10.mass, y = log10.hra)) +  
  ggplot2::geom_point()
```



Create a similar scatterplot showing the relationship between the raw values rather than the log values.

4. Colorize scatterplot points by class:

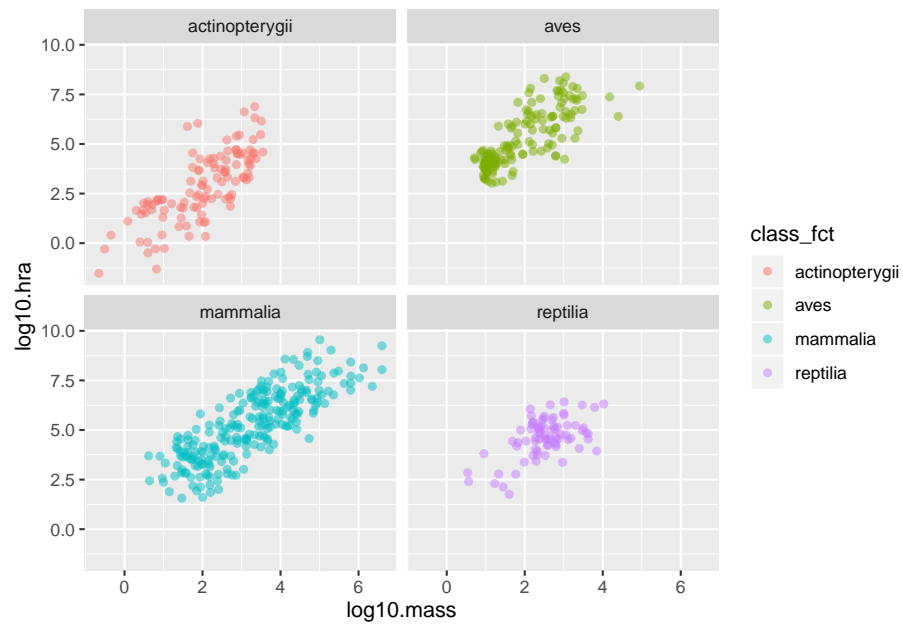
```
hrra %>%  
  dplyr::mutate(class_fct = as.factor(class)) %>%  
  ggplot2::ggplot(mapping = aes(x = log10.mass, y = log10.hrra, color = class_fct)) +  
  ggplot2::geom_point(alpha = 0.5)
```



*Group by order and experiment with different alpha values.*

5. Create a faceted plot:

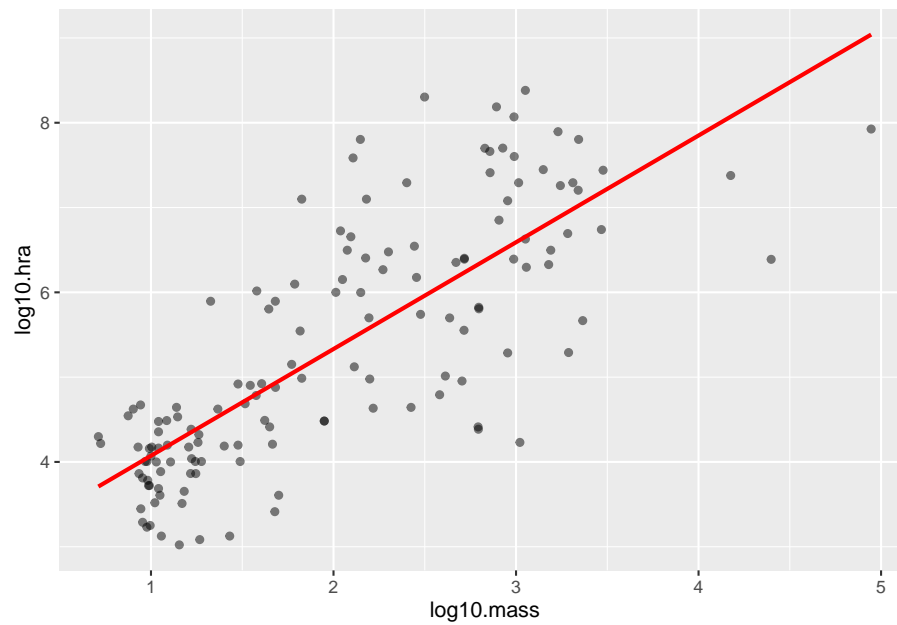
```
hrra %>%
  dplyr::mutate(class_fct = as.factor(class)) %>%
  ggplot2::ggplot(mapping = aes(x = log10.mass, y = log10.hra, color = class_fct)) +
  ggplot2::geom_point(alpha = 0.5) +
  ggplot2::facet_wrap(vars(class_fct))
```



Create a plot faceted by order for just the reptiles.

6. Fit a linear regression to the logarithmic data for birds:

```
hra %>%
  dplyr::filter(class == "aves") %>%
  ggplot2::ggplot(mapping = aes(x = log10.mass, y = log10.hra)) +
  ggplot2::geom_point(alpha = 0.5) +
  ggplot2::geom_smooth(method = lm, color = 'red', se = FALSE)
```



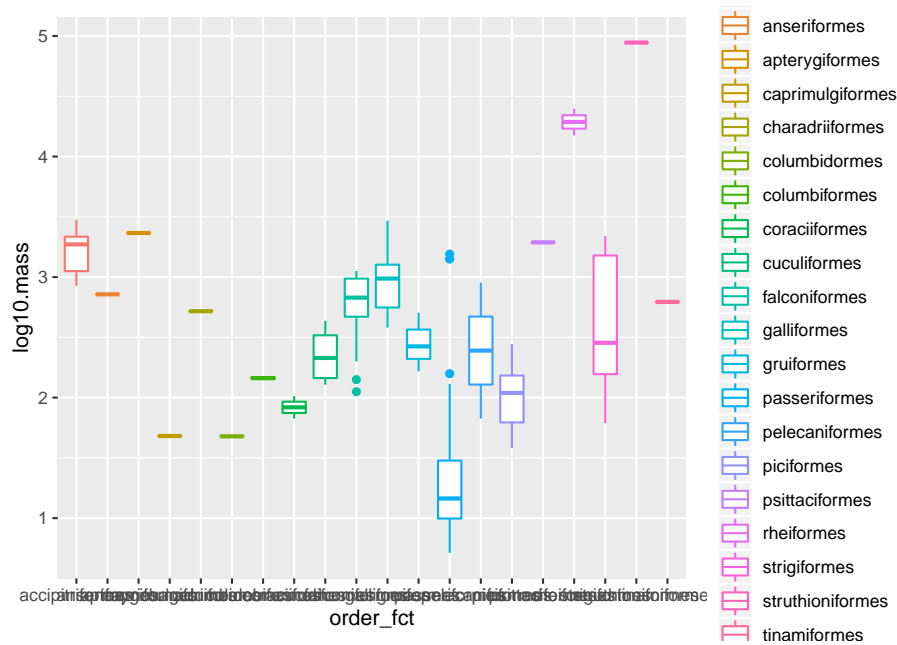
*Fit a line to the raw data for birds rather than the logarithmic data.*

7. Create a violin plot of mass by order for birds:

```
hra %>%  
  dplyr::filter(class == "aves") %>%  
  dplyr::mutate(order_fct = as.factor(order)) %>%  
  ggplot2::ggplot(mapping = aes(x = order_fct, y = log10.mass, color = order_fct)) +  
  ggplot2::geom_violin()
```



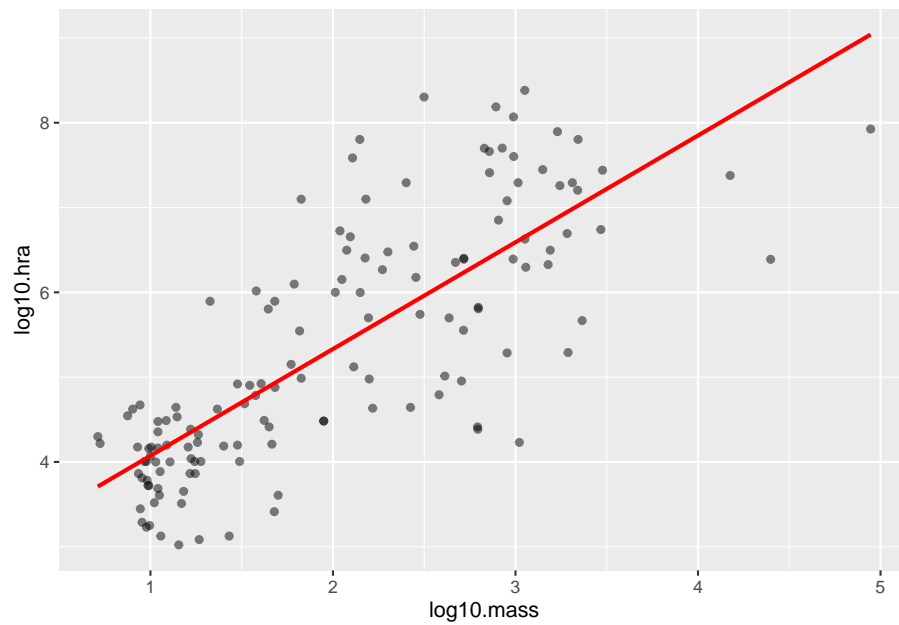




*Fix the labels and remove orders that only contain one species.*

9. Save the linear regression plot for birds as a PNG:

```
hra %>%
  dplyr::filter(class == "aves") %>%
  ggplot2::ggplot(mapping = aes(x = log10.mass, y = log10.hra)) +
  ggplot2::geom_point(alpha = 0.5) +
  ggplot2::geom_smooth(method = lm, color = 'red', se = FALSE)
```



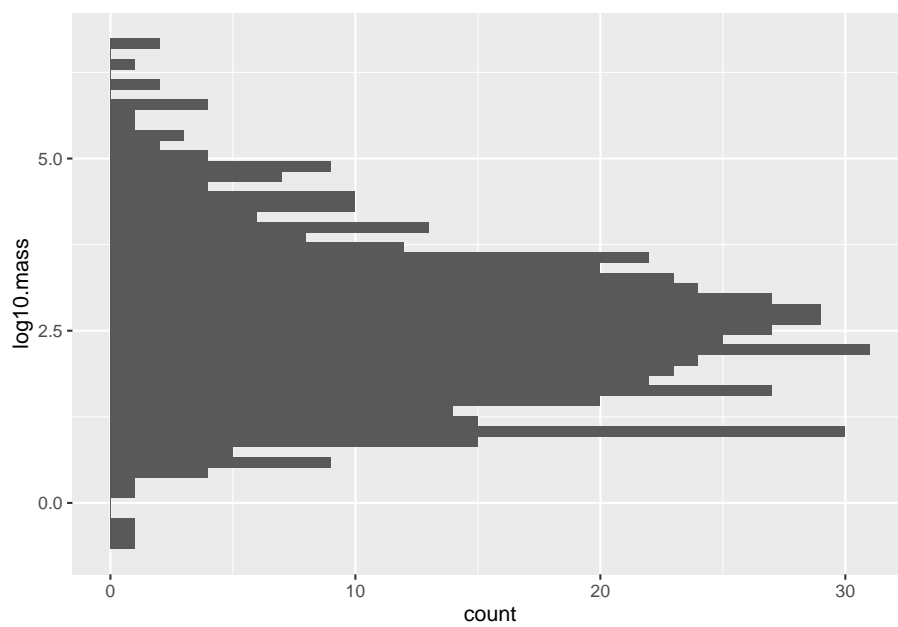
```
ggsave(here::here("birds.png"))
```

```
## Saving 6.5 x 4.5 in image
```

*Save the plot as SVG scaled to be 8cm wide.*

10. Create a horizontal histogram with 50 bins:

```
ggplot2::ggplot(hra, mapping = aes(x = log10.mass)) +  
  ggplot2::geom_histogram(bins = 50) +  
  ggplot2::coord_flip()
```



Use `stat_summary` to summarize the relationship between mass and home range area by class.



## Chapter 5

# Reproducibility

### 5.1 Questions {r-reproducibility-questions}

- How do I organize my data analysis projects, my files, and my folders?
- How can I make sure that my code is understandable to others?
- How do I ensure my analyses are reproducible, by others and myself in the (not-too-distant) future?

### 5.2 Objectives

1. To use a standard file and folder structure when starting a new project for better project management and organization.
2. To use RStudio to manage projects in a consistent and structured way.
3. To structure your data analysis projects to be more reproducible.
4. To create reproducible documents interwoven with R code that can be easily updated by changing the code or data.

### 5.3 Introduction

#### 5.3.1 What is reproducibility and why try to do it?

There are several key cornerstones for doing rigorous and sound scientific research, two of which are reproducibility and replicability (Patil et al., 2016). Replicability is when a study is repeated by other independent research

groups. Reproducibility is when, given the same data and the same analytical/computational steps, a scientific result can be verified. Both of these concepts are surprisingly difficult to achieve.

This course is about data analysis, so we'll be focusing solely on reproducibility rather than replicability. At present, there is little effort in science for having research be reproducible, likely due in many ways to a lack of training and awareness. Being reproducible isn't just about doing better science, it can also:

1. Make you much more efficient and productive, as you spend less time between coding and putting your results in the document.
2. Make you more confident in your results, since what you report and show as figures or tables will be exactly what you get from your analysis. No copying and pasting required!

There are many aspects to reproducibility, such as:

- Organized files and folder, preferably based on a standard or conventional structure.
- Understandable and readable code that is documented and descriptive.
- Results from analyses are identical to results presented in scientific output (e.g. article, poster, slides).
- Results from analyses are identical when code is executed on other machines (results aren't dependent on one computer).
- Explicit description or instruction on the order that code and scripts need to be executed.

We'll cover the first three items in this course.

## 5.4 Project organization

First off, what exactly does “project” mean? That depends a bit on the group, individual, or situation, but for our purposes, a “project” is anything related to one or more completed scientific “products” (e.g. poster, slides, manuscript, package, teaching material) related to a specific question or goal. This could be “one manuscript publication and associated conference presentations” per project. Confining a project to one “scientific output” facilitates keeping the project reproducible is kept reproducible, all files will relate to that “output”, and can be easily archived once the manuscript has been published. However, this definition could be different depending on your own situation and goals.

The ability to read, understand, modify, and write simple pieces of code is an essential skill for modern data analysis tasks and projects. Here we introduce you to some of the best practices one should have while writing their code. , many of which were taken from published “best practices” articles (Noble, 2009; Taschuk and Wilson, 2017; Wilson et al., 2017).

- Organise all R scripts and files in a single parent directory using a common and consistent folder and file structure.
- Use version control to track changes to files.
- Make raw data “read-only” (don’t edit it directly) and use code to show what was done.
- Write and describe code for people to read by being descriptive and using a style guide.
- Think of code as part of your research product: Write for an audience or other reader.
- Create functions to avoid repetition.
- Whenever possible, use code to create output (figures, tables) rather than manually creating or editing them.

Managing your projects in a reproducible fashion doesn’t just make your science reproducible, it also makes your life easier! RStudio is here to help us with that by using R Projects. RStudio projects make it straightforward to divide your work into multiple contexts, each with their own working directory, workspace, history, and source documents.

It is strongly recommended that you store *all* the necessary files that will be used in your code in the **same parent directory**. You can then use relative file paths to access them (we’ll talk about file paths below). This makes the directory and R Project a “product” or “bundle/package”. Like a tiny machine, that needs to have all its parts in the same place.

### 5.4.1 Creating your first project

There are many ways one could organise a project folder. We’ll set up a project directory folder using the `prodigenr` package:

```
# prodigenr::setup_project("ProjectName")
prodigenr::setup_project("learningr")
```

(You can also create a new project in RStudio by using “File -> New Project -> Scientific Analysis Project using prodigenr”.)

When we use the `::` colon here, we are telling R “use `setup_project` function from the `prodigenr` package”. This function will then create the following folders and files:

```
learningr
├── R
│   ├── README.md
│   ├── fetch_data.R
│   └── setup.R
└── data
    └── README.md
```

```
doc
  README.md
.gitignore
DESCRIPTION
learningr.Rproj
README.md
```

This forces a specific, and consistent, folder structure to all your work. Think of this like the “introduction”, “methods”, “results”, and “discussion” sections of your paper. Each project is then like a single manuscript or report, that contains everything relevant to that specific project. There is a lot of power in something as simple as a consistent structure. Projects are used to make life easier. Once a project is opened within RStudio the following actions are taken:

- A new R session (process) is started.
- The current working directory is set to the project directory.
- RStudio project options are loaded.

The README in each folder explains a bit about what should be placed there. But briefly:

1. Documents like manuscripts, abstracts, and exploration type documents should be put in the `doc/` directory (including R Markdown files which we will cover later).
2. Data, raw data, and metadata should be in either the `data/` directory or in `data-raw/` for the raw data.
3. All R files and code should be in the `R/` directory.
4. Name all new files to reflect their content or function. Follow the tidyverse style guide for file naming.

Note the `DESCRIPTION` file. This is used as metadata about the project and is useful when working on R projects. For any project, it is **highly recommended** to use version control. We’ll be covering version control in more detail later in the course.

### 5.4.2 Exercise: Better file naming

Look at the list of file names below. Which file names are good names and which shouldn’t you use?

```
fit models.R
fit-models.R
foo.r
stuff.r
get_data.R
Manuscript version 10.docx
manuscript.docx
new version of analysis.R
```



```
trying.something.here.R  
plotting-regression.R  
utility_functions.R  
code.R
```

### 5.4.3 Should you keep your data under version control?

We have a `data/` folder for a reason. But you might not want to keep the data under version control, for several reasons:

1. It's a large dataset (tens or more Mb in file size)
2. There are sensitive and/or personally-identifying information in the data

As a rule of thumb, if you can send the data by an email attachment, you could probably put it into Git. Unless there is sensitive or personal data, then don't. If it isn't kept under version control, make sure you include a reference to how or where you got the data, either as an R script showing the code you used to import/clean/download it or described in the `README.md` file.

## 5.5 Reusability

Part of reproducibility is also making sure your scripts and file organization is “reusable” meaning that others (or yourself) can run it again. So, for instance, making sure to use “relative file paths” compared to “absolute file paths” (we'll cover these in a bit). Or indicating which other R packages your code depends on. So here we'll cover how to make sure your scripts and project files are reusable.

### 5.5.1 Keeping a clean slate

When you finish writing your R code for the day and close the session, you probably will be asked about saving your session. What this does is everything kept in the environment (e.g. all objects, functions, or datasets you created and used during the session) get saved to an `.RData` file. Then, the next time you open up your R session, R will see this `.RData` file and load everything in that file. Everything you did previously will be loaded into your environment. This seems like a good thing... but it's not. Imagine eating your dinner on a really dirty plate... that's not pleasant right? Loading a previous session is like that dirty plate.

So, to make sure you always use a clean slate, we'll run a handy function from the `usethis` package to stop R from saving and loading this `.RData` file, ensuring you have a clean working environment. You only need to run this function once, as it will set the appropriate RStudio settings for you.

```
usethis::use_blank_slate()
```

We'll use the `usethis` package more throughout this chapter and others, as it provides several very useful functions when working with projects.

### 5.5.2 Packages, data, and file paths

A major strength of R is in its ability for others to easily create packages that simplify doing complex tasks (e.g. creating figures with the `ggplot2` package) and for anyone to easily install and use that package<sup>1</sup>. You load a package by writing:

```
library(tidyverse)
```

Working with multiple R scripts and files, it quickly gets tedious to always write out each library function at the top of each script. A better way of managing this is to create a new file, keep all package loading code in that file, and sourcing that file in each R script. So, to create a new R file in the `R/` folder, we'll use this `use_r()` function from the `usethis` package:

```
usethis::use_r("package-loading")
```

This creates a file called `package-loading.R` in the `R/` folder. In the `package-loading.R` file, add this code to it.

```
library(tidyverse)
```

Then in other R scripts in the `R/`, include this code at the top the script:

```
source(here::here("R/package-loading.R"))
```

The `here` package uses a function called `here()` that makes it easier to manage file paths. What is a file path and why is this necessary? A file path is the list of folders a file is found in. For instance, your resume may be found in `/Users/Documents/personal_things/resume.docx`. The problem with file paths in R is that when you run a script interactively (e.g. what we do in class and normally), the file path is located at the Project level (where the `.Rproj` file is found). You can see the file path by looking at the top of the "Console".

But! When you `source()` an R script, it may likely run *in the folder it is saved in*, e.g. in the `R/` folder. So your file path `R/packages-loading.R` won't work because there isn't a folder called `R` in the `R/` folder. Often people use the function `setwd()`, but this is *never* a good idea since using it makes your script

---

<sup>1</sup>You may encounter some who say you shouldn't rely on packages and to only use base R functions. However, this is seriously bad advice since the ecosystem of R packages can greatly simplify your life doing data analysis. Plus, packages greatly expand and enhance the capability of R, so make use of packages! If someone invents a wheel, why wouldn't you use it?

*runnable only on your computer...* which makes it no longer reproducible. We use the `here()` function to tell R to go to the project root (where the `.Rproj` file is found) and then use the file path from that point. This simple function can make your work more reproducible and easier for you to use later on.

We also use the `here()` function when we import a dataset or save a dataset. So, let's load in the NYC Dog License dataset. First, save the CSV in the `data/` folder. Then create a new R file:

```
usethis::use_r("load-data")
```

And write these lines in the file:

```
source(here::here("R/package-loading.R"))
dog_license <- read_csv(here::here("data/nyc-dog-licenses.csv.gz"))
head(dog_license)
```

That is how we will load data in from now on.

Here are a few other tips for keeping your code reusable:

- When encountering a difficult problem, try to find R packages or functions that do your problem for you<sup>2</sup>.
- Split up your analyses steps into individual files (e.g. “model” file, “plot” file). Then `source` those files as needed or save the output in `data/` to use it in other files.
- Try not to have R scripts be too long (e.g. more than 500 lines of code). Keep script sizes as small as necessary and as specific as possible (have a single purpose). A script should have an end goal.

## 5.6 Readability

There are two reasons we write code: to instruct the computer to do something and to record the steps we took to get a particular result for us or others to understand. For computers, *how* or *what* you write doesn't matter, as long as the code is correct. Computers don't need to *understand* the code. But humans do need to understand it. We need clear language and explicit meaning in order to understand what is going on. Humans write code, humans read code, and humans must maintain it and fix any errors. So, what you write and how you write it is extremely important.

Like natural human languages, R has a relaxed approach to how R code is written. This has some nice advantages, but also some major disadvantages, notably that writing styles can be quite different across the world or even within

---

<sup>2</sup>You may hear some people say “oh, don't bother with R packages, do everything in base R”... don't listen to them. Do you build a computer from scratch everytime you want to do any type of work? Or a car when you want to drive somewhere? No, you don't. Make use of other people's hard work to make tools that simplify your life.

one's own code. So, it's important to stick to some guidelines, for instance, as laid out by the tidyverse style guide. Some other tips include:

- Write your code assuming other people will be reading it.
- Stick to a *style guide*. (We're repeating this because it's really important!)
- Use full and descriptive words when typing and creating objects.
- Use white space to separate concepts (empty lines between them, use spaces, and/or tabs).
- Use RStudio R Script Sections ("Code->Insert Section" or Ctrl-Shift-R) to separate content in scripts.

Even though R doesn't care about naming, spacing, and indenting, it really matters how your code looks. Coding is just like writing. Even though you may go through a brainstorming note-taking stage of writing, you eventually need to write correctly so others can understand, *and read*, what you are trying to say. Brainstorming and exploratory work is fine, but eventually you need to write code that will be legible. That's why using a style guide is really important.

### 5.6.1 Exercise: Make the code more readable

Using the style guide found in the link, try to make these code more readable. Edit the code so it follows the correct style and so it is easier to understand and read. You don't need to understand what the code does, just follow the guide.

```
# Object names
DayOne
dayone
T <- FALSE
c <- 9
mean <- function(x) sum(x)

# Spacing
x[,1]
x[ ,1]
x[ , 1]
mean (x, na.rm = TRUE)
mean( x, na.rm = TRUE )
function (x) {}
function(x){}
height<-feet*12+inches
mean(x, na.rm=10)
sqrt(x ^ 2 + y ^ 2)
df $ z
x <- 1 : 10

# Indenting
```

```
if (y < 0 && debug)
  message("Y is negative")
```

FIXME: The below “details” will need to be dealt with since PDF doesn’t allow this TODO: Maybe move to a solutions section at the end of chapter?

Click for a possible solution

The old code is in comments and the better code is below it.

```
# Object names

# Should be camel case
# DayOne
day_one
# dayone
day_one

# Should not over write existing function names
# T = TRUE, so don't name anything T
# T <- FALSE
false <- FALSE
# c is a function name already. Plus c is not descriptive
# c <- 9
number_value <- 9
# mean is a function, plus does not describe the function which is sum
# mean <- function(x) sum(x)
sum_vector <- function(x) sum(x)

# Spacing
# Commas should be in correct place
# x[,1]
# x[ ,1]
# x[ , 1]
x[, 1]
# Spaces should be in correct place
# mean (x, na.rm = TRUE)
# mean( x, na.rm = TRUE )
mean(x, na.rm = TRUE)
# function (x) {}
# function(x){}
function(x) {}
# height<-feet*12+inches
height <- feet * 12 + inches
# mean(x, na.rm=10)
mean(x, na.rm = 10)
# sqrt(x ^ 2 + y ^ 2)
```

```

sqrt(x^2 + y^2)
# df $ z
df$z
# x <- 1 : 10
x <- 1:10

# Indenting should be done after if, for, else functions
# if (y < 0 && debug)
# message("Y is negative")
if (y < 0 && debug)
  message("Y is negative")

```

### 5.6.2 Automatic styling with styler

You may have organised the exercise by hand, but it's possible to do it automatically. The tidyverse style guide has been implemented into the `styler` package to automate the process of following the guide by directly re-styling selected code. The `styler` snippets can be found in the `Addins` function in the RStudio “Addins” menu after you have installed it.

RStudio also has its own automatic styling ability, through the menu item “Code -> Reformat Code” (or `Ctrl-Shift-A`). Try both methods of styling on the exercise code above. There are slight differences in how each method works and they both aren't always perfect. For now, let's stick with using the `styler` package. The `styler` functions work on R code within both `.R` script files as well as R code within `.Rmd` documents, which we will cover later in this lesson.

There are several `styler` RStudio addins, but we'll focus on the two:

- “Style selection”: Highlight text and click this button to reformat the code.
- “Style active file”: Code in the `.R` or `.Rmd` file you have open and visible in RStudio will be reformatted when you click this button.

There are two other `styler` functions that are also useful:

- `styler::style_file("path/to/filename")`: Styles the whole file as indicated by the file path in the first argument. Can be either an `.R` or `.Rmd` file.
- `styler::style_dir("directoryname")`: Styles all files in the indicated directory in the first argument.

Let's try the `styler::style_file()` function out. Inside a file called `non-styled-code.R`, it has:

```

# Spacing
x[,1]

```

```

mean (x, na.rm = TRUE)
function (x) {}
height<-feet*12+inches
sqrt(x ^ 2 + y ^ 2)
df $ z
x <- 1 : 10

# Indenting
if (y < 0 && debug)
message("Y is negative")

```

Then we run:

```
styler::style_file("testing-styler.R")
```

```

Styling 1 files:
  testing-styler.R

```

```

Status Count Legend
    0   File unchanged.
    1   File changed.
    0   Styling threw an error.

```

Please review the changes carefully!

Which changes the file to be styled!

```

# Spacing
x[, 1]
mean(x, na.rm = TRUE)
function(x) {}
height <- feet * 12 + inches
sqrt(x^2 + y^2)
df$ z
x <- 1:10

# Indenting
if (y < 0 && debug) {
  message("Y is negative")
}

```

This is more or less everything that the styler package does.

### 5.6.3 Exercise: Use styler to fix code formatting

Use the styler package function on the code from the previous exercise by either running `styler::style_file()` or with the "Style selection" addin when

highlighting the code.

## 5.7 Integrating text, code, and results

The most obvious demonstration of reproducibility is when the results obtained from executing the analysis code (by an independent entity) are identical to the results presented in the scientific output such as in an article. When there is agreement between these two, reproducibility has been more or less achieved. In R, there are tools available to *completely* ensure that this happens by directly inserting the results from the code *into the scientific output*. This is done by using R Markdown, which interweaves R code with text. So instead of, for example, manually inserting a figure, you write R code within the document to insert the figure for you! Using R Markdown can save so much time and get your work that much closer to being reproducible.

There are many other advantages to using R Markdown. From the single R Markdown format you can use it to create manuscripts, posters, slides, websites, books, and many more from simply using R Markdown. In fact, this book was written using R Markdown. As a bonus, switching between citation formats or Word templates for different journals is easier than doing it with Word.

### 5.7.1 Markdown

R Markdown uses, well, Markdown as the format to convert to multiple document types. Fun fact: This website is built based on Markdown! While there are many “flavours” of Markdown that have been developed over the years, R Markdown uses the Pandoc version. Pandoc is a combination of pan which is Latin for “all” and doc which means document.

Markdown is a “markup language” meaning that special characters mean certain things, which we will cover below.

To format text, such as to bold or make a list, you use the special characters. You write Markdown as plain text (like R code), so you don’t need any special software (like you do with Word documents). Most features needed for writing a scientific document are available in Markdown, but not all. *Tip:* Try to fit your writing and document creation around what Markdown can achieve, rather than force or fight Markdown to do something it wasn’t designed to do.

All right, let’s create and save an R Markdown file. In RStudio, go to **File -> New File -> R Markdown**. A dialog box will pop up. In the “Title” section, type in **Reproducible documents** and in the “Author” section type in your name. Choose the HTML output format. Save this file as **learning-rmarkdown.Rmd** in the **doc/** folder.



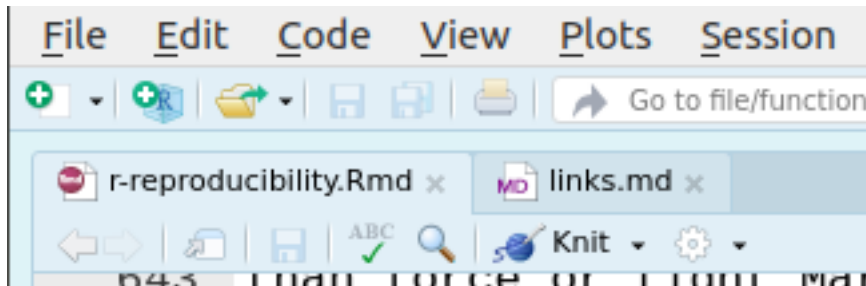


Figure 5.1: “Knit” button.

Inside the file, there is a bunch of text that shows some basic formatting you can use for writing Markdown. For now, delete everything *except* the top part of the file (the part surrounded by `---`). This part is called the YAML header, which we will cover more a bit later. Try converting the file to HTML by hitting the “Knit” button at the top or by typing out `Ctrl-Shift-K`.

#### 5.7.1.1 Headers

Creating headers (like chapters or sections) is indicated by using one or more `#` at the beginning of a line, prefixing some text:

```
# Header level 1

Paragraph.

## Header level 2

Paragraph.

### Header level 3

Paragraph.
```

This creates the section headers as seen directly above (“Headers”) or below (“Text formatting”). The header text *must* be on one line, otherwise the next line is interpreted as paragraph text.

See how it looks by “Knitting” the document (“Knit” button or `Ctrl-Shift-K`).

#### 5.7.1.2 Text formatting

To format text individually, surround the text with the special characters, as shown here:

- **`**bold**`** gives **bold**.
- *`*italics*`* gives *italics*.
- `superscript` gives super<sup>script</sup>.
- `subscript` gives sub<sub>script</sub>.

What if you want to use the special character as simple text? Prefix it with an `\`, so `\*` becomes `*`, `\^` becomes `^`, and `\~` becomes `~`.

### 5.7.1.3 Lists

To create an unnumbered list, do:

```
- item 1
- item 2
- item 3
```

which gives...

- item 1
- item 2
- item 3

Notice the empty lines above and below the list, those are important. To create a numbered list, do:

```
1. item 1
2. item 2
3. item 3
```

which gives...

1. item 1
2. item 2
3. item 3

See how it looks by “Knitting” the document (“Knit” button or `Ctrl-Shift-K`).

### 5.7.1.4 Blockquotes

Sometimes (probably not too commonly), you may need to quote someone by using “blockquotes”. To do that, do:

```
> Blockquote
```

which gives...

Blockquote

Blockquotes can be as many lines as you want. To stop a paragraph from being in the blockquote, separate the text with an empty line:

```
> Blockquote paragraph
```

```
Regular paragraph
```

#### 5.7.1.5 Adding footnotes

Footnotes can be added by using `[^some-text-label]`, such as:

```
Footnote[^1]
```

```
[^1]: Footnote content.
```

which gives...

Footnote<sup>3</sup>

So you can write some text, add some footnotes within, and include the footnote content right below the paragraph:

```
Paragraph text[^1], with some more text[^reference].
```

```
[^1]: This is the first footnote.
```

```
[^reference]: This is the next footnote.
```

```
More paragraphs.
```

Notice the empty lines in between. The footnote should also be on one line, though it isn't strictly necessary. See how it looks by “Knitting” the document (“Knit” button or `Ctrl-Shift-K`).

#### 5.7.1.6 Inserting pictures, images, or figures

You can include externally created (i.e. not by an R code chunk, discussed later on) png, jpeg, or pdf image file by adding (:

```
![Image caption here.](path/to/image/file.png)
```

So something like this:

```
![Steps to being more reproducible. Source DOI: 10.1038/d41586-018-05990-5](figures/r-reproducibi
```

which gives...

*Tip:* Can also include links to images from the Internet, as a URL link.

---

<sup>3</sup>Footnote content.

## Getting reproducible

There's no one-size-fits-all solution for computational reproducibility. But these practices can help.

- **Use code.** Instead of pointing and clicking, use programming languages to download, filter, process and output your data, and command-line scripts to document how those tools are executed.

- **Go open-source.** Code transparency is key to reproducibility, so use open-source tools whenever possible. "If you give me a black box with no source code and it just gives me numbers, as far as I am concerned, it's a random-number generator," says mathematician Les Hatton of Kingston University in London.

- **Track your versions.** Using version-control software such as Git and GitHub, researchers can document and share the precise version of the tools that they use, and retrieve specific versions as necessary.

- **Document your analyses.** Use computational notebooks such as Jupyter to interleave code, results and explanatory text in a single file.

- **Archive your data.** Freeze data sets at key points — when submitting an article for publication, for example — with archiving services such as Zenodo, Figshare or the Open Science Framework.

- **Replicate your environment.** Software 'containers', such as Docker, bundle code, data and a computing environment into a single package; by unboxing it, users can recreate the developer's system. ReproZip, developed in the lab of New York University computer scientist Juliana Freire, simplifies container creation by watching program execution to identify its requirements. The commercial service Code Ocean and an open-source alternative, Binder, enable researchers to create and share executable Docker containers that users can explore in a web browser.

- **Automate.** Automation provides reproducibility without users really having to think about it, says bioinformatician Casey Greene at the University of Pennsylvania in Philadelphia. Continuous integration services such as Travis CI can automate quality-control checks, for instance, and the Galaxy biocomputing environment automatically logs details of the jobs it runs.

- **Get help.** Resources abound for interested researchers; see [practicereproducibleresearch.org](http://practicereproducibleresearch.org), for instance, or find a Software Carpentry workshop near you to learn basic computing skills.

Figure 5.2: Steps to being more reproducible. Source DOI: 10.1038/d41586-018-05990-5



Figure 5.3: Caption.

If you want to modify the width or sizing, append something like `{width=##%}` to the end of the image insertion:

```
! [Caption.] (figures/r-reproducibility/code-sharing-steps.png) {width=50%}
```

which gives...

#### 5.7.1.7 Adding links to websites

And a link can be linked in the following format:

```
[Link to GitHub] (https://github.com).
```

gives...

Link to GitHub.

The above form is great for one time use, but what if you want to use the same link again later on? Use this form then:

Use multiple `[links]` in your document. The same `[links]` can be used again.

```
[links]: https://github.com
```

which gives...

Use multiple links in your document. The same links can be used again.

### 5.7.1.8 Inserting (simple) tables

You can insert tables with Markdown too. We wouldn't recommend doing it for complicated tables though, as it can get tedious fast! (A recommended approach for more complex or bigger tables is to make the table contents as a data frame in R first and then use the `knitr::kable()` function to create the table, as we'll cover in the R Markdown section below). You can even include Markdown text formatting inside the table:

```
|   | Fun | Serious |
|:--|--:|:--:|
| **Happy** | 1234 | 5678 |
| **Sad** | 123 | 456 |
```

which gives...

	Fun	Serious
<b>Happy</b>	1234	5678
<b>Sad</b>	123	456

The `|:--:|` or `|:--|` tell the table to right-align or left-align, respectively, the values in the table column. Center-align is `|:--:|`. See how it looks by “Knitting” the document (“Knit” button or `Ctrl-Shift-K`).

### 5.7.1.9 Exercise: Try to re-create a document using Markdown

1. Open this link. This is a HTML file that has been created by using Markdown formatting.
2. Create a new R Markdown file, with output type “HTML”, and save it as `mimic-html-file.Rmd`.
3. Delete all the automatically added text except the top part Inside the R Markdown file
4. Write text using Markdown formatting so that you can create a `html_document` that looks exactly like the linked file.
5. Knit the R Markdown document. Confirm that your version looks the same as the above version.

## 5.7.2 R Markdown

R Markdown is an extension of Markdown that weaves together R code with Markdown formatted text all together in a single document. Output from the R code gets inserted directly into the document for a seamless integration of document writing and analysis reporting.

### 5.7.2.1 YAML header/metadata

Most Markdown documents (especially for R Markdown) include YAML metadata at the top of the document, surrounded by `---`. YAML is a data format, like CSV, that contains the metadata and various options that R Markdown uses for the entire document. Data in YAML is stored in the form `variable: value`. For instance, `title` or `author` is paired with their respective “values” and then used by Markdown when creating the document. Other options are also included here, such as what the converted output document should be, such as Word. There are many more output formats to choose from (e.g. slides, websites, books). The YAML header looks something like this:

```
---
title: "Document title"
author: Your Name
output: html_document
---
```

Here, there are three variable-value pairings: `title`, `author`, and `output`. In the `output` variable, the R Markdown function `html_document` is given so that the output document format is converted to HTML. There are also `word_document` and `pdf_document` settings. For now, we’ll focus on the `html_document` output. Usually when you create the R Markdown file, this YAML header gets added automatically. There are additional options you can set in the output field, which we will cover later on.

### 5.7.2.2 Using R code chunks

R Markdown’s primary function is to allow combining text and R code in the same document, which is incredibly powerful and useful! All R code chunks have the appearance:

```
```{r chunk-label-name, chunk.option="...", chunk.option=...}
...R code...
```
```

Notice that chunk options need to be on the same line (one single line). Any R code in the code chunk gets evaluated and any output gets inserted into the document. So if the code prints to the console, it will get inserted into the document.

The `r` tells R Markdown to run R on the code chunk, while the chunk label differentiates the chunk from other chunks. Using a chunk label also helps navigate a document when you use RStudio’s “Document outline” (`Code -> Show Document Outline`). It also does a few other things.

**Note:** Standard practice is that code chunk labels should be named without `_`, spaces, or `.` and instead should be one word or be separated by `-`.

Let's load in some packages and data:

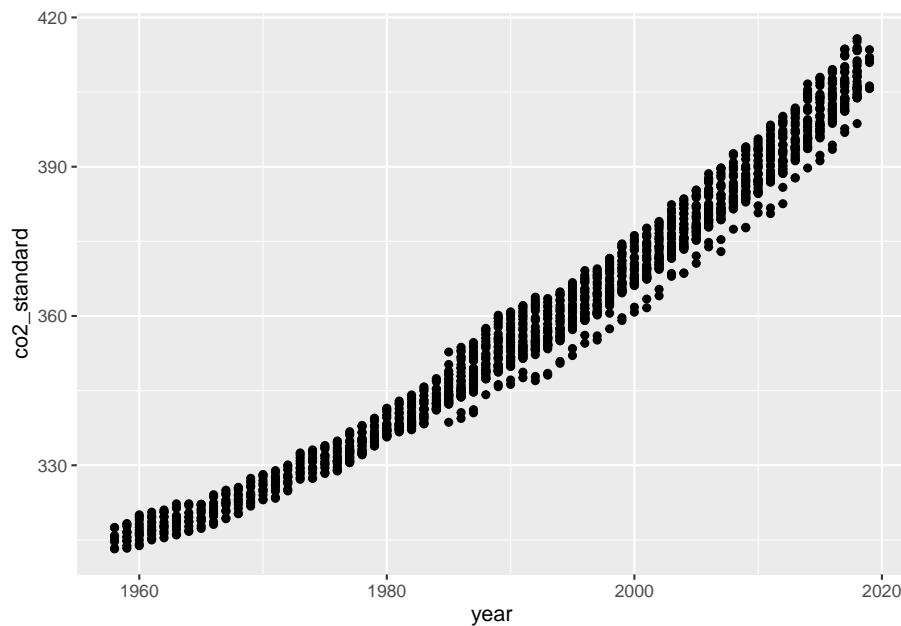
```
```{r setup}
library(tidyverse)
co2_data <- read_csv(here::here("data/co2.csv")) %>%
  filter(co2_standard > 0)
```
```

When running your code chunks interactively as you develop and write your document, R Markdown will look for a code chunk labeled `setup` to run first (for instance, to load all packages used in a document). Hence we name this chunk “setup”.

### 5.7.2.3 Inserting figures

One of the most obvious benefits to using R Markdown is to automatically insert a plot. To do that we do:

```
```{r co2-time-plot}
ggplot(co2_data, aes(x = year, y = co2_standard)) +
  geom_point()
```
```



What if we want to include the plot but not the code? Easy! Set the chunk option `echo` to `FALSE`:



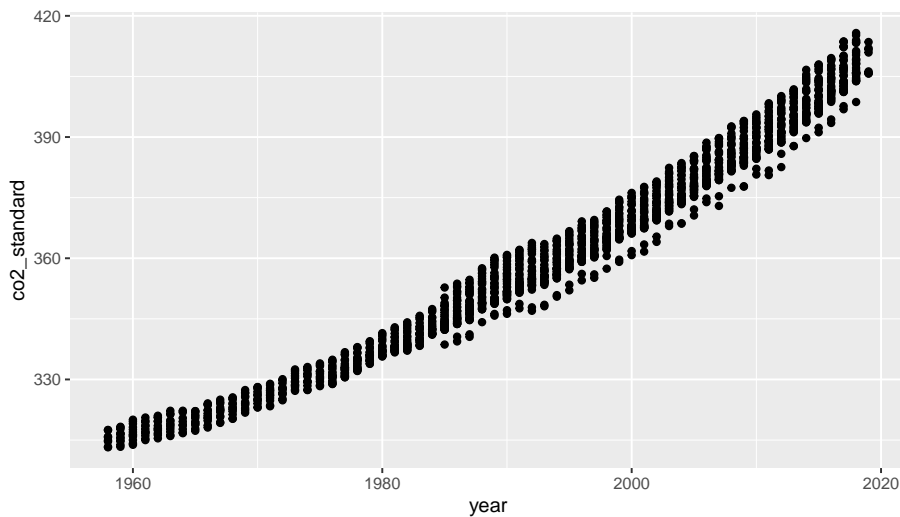


Figure 5.4: Add your figure title here.

```
```{r co2-time-plot, echo=FALSE}
ggplot(co2_data, aes(x = year, y = co2_standard)) +
  geom_point()
```
```

Or what if you don't want R to run the code chunk, but still show the code? Set the chunk option `eval` (for evaluate) to `FALSE`:

```
```{r co2-time-plot, eval=FALSE}
ggplot(co2_data, aes(x = year, y = co2_standard)) +
  geom_point()
```
```

Since we have a figure, we can change some width, height, and alignment options, as well as add a caption with the `fig.width`, `fig.height`, `fig.align`, and `fig.cap` (respectively):

```
```{r co2-time-plot, fig.cap="Add your figure title here.", fig.height=4, fig.width=7}
ggplot(co2_data, aes(x = year, y = co2_standard)) +
  geom_point()
```

ggplot(co2_data, aes(x = year, y = co2_standard)) +
  geom_point()
```

What if we want to first run some code before running the plot, but don't want to include the output of the other code? Make sure that the output doesn't

“print”. So letting R evaluate both dataframe and plot, the code chunk will output and insert both into the knitted document.

```
```{r co2-time-plot-after-2000-1, fig.cap="Add your figure title here.", fig.height=4,
co2_after_2000 <- co2_data %>%
  filter(year >= 2000)

# Print dataframe
co2_after_2000

# Print plot
ggplot(co2_after_2000, aes(x = year, y = co2_standard)) +
  geom_point()
```
```

```
## # A tibble: 680 x 5
##   year month date_numeric co2_standard station
##   <dbl> <dbl>      <dbl>      <dbl> <chr>
## 1 2000     1      2000.        369. Mauna Loa, Hawaii, USA
## 2 2000     2      2000.        369. Mauna Loa, Hawaii, USA
## 3 2000     3      2000.        371. Mauna Loa, Hawaii, USA
## 4 2000     4      2000.        372. Mauna Loa, Hawaii, USA
## 5 2000     5      2000.        372. Mauna Loa, Hawaii, USA
## 6 2000     6      2000.        372. Mauna Loa, Hawaii, USA
## 7 2000     7      2001.        370. Mauna Loa, Hawaii, USA
## 8 2000     8      2001.        368. Mauna Loa, Hawaii, USA
## 9 2000     9      2001.        367. Mauna Loa, Hawaii, USA
## 10 2000    10      2001.        367. Mauna Loa, Hawaii, USA
## # ... with 670 more rows
```

Compare to this next code chunk, which will only output (“print”) the plot. See how we don’t include code that would send anything to the console to be “printed”? Only things that get “printed” will be included in the R Markdown output document. Also notice how we renamed the chunk label? In R Markdown you can’t have duplicate code chunk labels.

```
```{r co2-time-plot-after-2000-2, fig.cap="Add your figure title here.", fig.height=4,
co2_after_2000 <- co2_data %>%
  filter(year >= 2000)

# Print plot
ggplot(co2_after_2000, aes(x = year, y = co2_standard)) +
  geom_point()
```
```

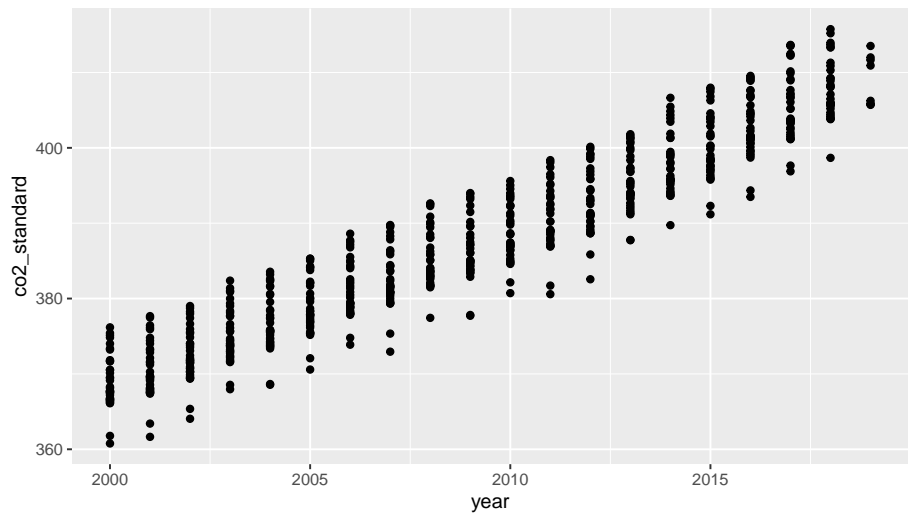


Figure 5.5: Add your figure title here.

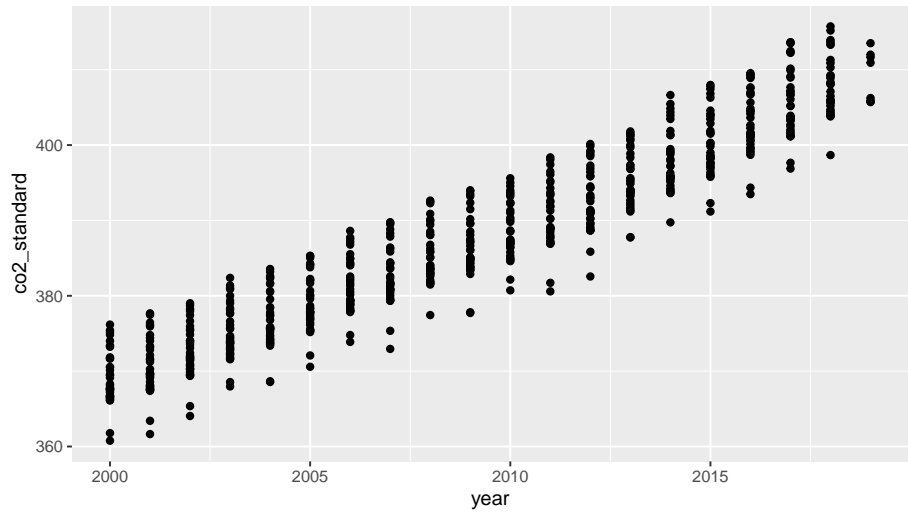


Figure 5.6: Add your figure title here.

### 5.7.2.4 Exercise: Add some figures to a R Markdown document

1. Create a new R Markdown file (“File -> New File -> R Markdown”), providing the title, author name (your name), and setting the output to HTML.
2. Save the file as `using-rmarkdown.Rmd` in the `doc/` folder.
3. Delete all text *except* the YAML header.
4. Create an R code chunk and call the label “setup”. Write code so the packages and data are loaded.
5. Create another code chunk and call it “plot-licenses-by-year”. Write R code to create a point plot (`geom_point()`) of the year on the x axis and number of licenses on the y axis.
6. Knit the document and see what the output looks like.
7. Change the theme of the plot to another builtin theme (*hint*: themes start with `theme_`).

### 5.7.2.5 Using R code chunks to insert tables

You can also create tables by using the `kable()` function from the `knitr` package. Let’s create a table of the mean CO<sub>2</sub> concentration over the years at each monthly period for each station.

```
```{r mean-co2-table}
co2_data %>%
  select(station, month, co2_standard) %>%
  group_by(station, month) %>%
  summarise(MeanCO2 = round(mean(co2_standard, na.rm = TRUE), 2)) %>%
  spread(station, MeanCO2) %>%
  knitr::kable(caption = "Table caption here.")
```
```

### 5.7.2.6 Inline R code

Often you might have results inside the text you are writing. Here you can include R code within the text so that the results are inserted directly into the document. It looks like:

The mean of CO<sub>2</sub> is ‘`r round(mean(co2_data$co2_standard, na.rm = TRUE), 2)`’.

Which gives...

The mean of CO<sub>2</sub> is `round(mean(co2_data$co2_standard, na.rm = TRUE), 2)`.

Keep in mind that inline R code can *only* insert a single number or character value, nothing more.

Table 5.2: Table caption here.

| month | Alert Station, NWT, Canada | Cape Grim, Tasmania, Australia | Mauna Loa, Hawaii, USA |
|-------|----------------------------|--------------------------------|------------------------|
| 1     | 378.54                     | 364.63                         | 354.61                 |
| 2     | 380.06                     | 364.65                         | 355.96                 |
| 3     | 380.75                     | 364.73                         | 356.11                 |
| 4     | 381.07                     | 364.89                         | 357.47                 |
| 5     | 380.36                     | 363.34                         | 356.57                 |
| 6     | 378.43                     | 363.70                         | 356.64                 |
| 7     | 371.33                     | 364.16                         | 354.50                 |
| 8     | 365.98                     | 364.60                         | 352.51                 |
| 9     | 365.11                     | 364.85                         | 350.85                 |
| 10    | 369.21                     | 364.90                         | 351.49                 |
| 11    | 374.77                     | 364.80                         | 352.24                 |
| 12    | 376.94                     | 364.68                         | 353.53                 |

### 5.7.2.7 Citing literature with R Markdown

No scientific writing is complete without being able to include references. If you want to insert a citation, use the Markdown key `[@Cone2016]`, which will look like (Conery, 2016). The text `Cone2016` is the key that the bibliography manager uses to identify a specific reference. Adding more references is done by separating by a `;`, so like `[@AuthorYear; @Author2Year; @Author3Year]`. The resulting citation reference will be inserted at the bottom of the document. To get the bibliography to work, you'll also need to add a line to the YAML header like this:

```
---
title: "My report"
author: "Me!"
bibliography: my_references.bib
---
```

The `my_references.bib` is a `.bib` file found in the same folder as the `.Rmd` file. So in our case, the `.bib` file is in the `doc/` folder. You can also use other bibliography manager files, such as EndNote. See this documentation for which bibliography managers can be used.

Since all references are appended to the bottom of the document, it's good to add a final "Reference" section header to the end of your file, like so:

```
# References
```

### 5.7.2.8 Making your report prettier

This part mostly applies to HTML-based and PDF<sup>4</sup> outputs, since programmatically modifying or setting templates in Word documents is rather difficult<sup>5</sup>. Changing broad features of a document can be done by setting the “theme” of the document. Add an option in the YAML metadata like:

```
---
title: "My report"
output:
  html_document:
    theme: sandstone
---
```

Check out the R Markdown [documentation][rmd-appearance] for more types of themes you can use for HTML documents, and advanced topics such as parameterized R Markdown documents. Most of the Bootswatch themes are available for use in R Markdown to HTML conversion.

Want to add a Table of Contents? Easy! Add `toc: true` to the YAML header:

```
---
title: "My report"
output:
  html_document:
    theme: sandstone
    toc: true
---
```

Adding a `toc` only works for PDF and HTML, but *not* Word documents.

### 5.7.2.9 Exercise: Add a summary table, inline results, and a prettier theme

1. Use the R Markdown file from the previous exercise (`using-rmarkdown.Rmd`).
2. Create three new header 1 `#` sections: Objective, Results, Conclusion.
3. Write in the “Objective” section an idea you have about the Dog License dataset. It can be as simple as “How many dogs are there in New Year City?”. Include an *italics* in this sentence.
4. Create three new code chunks in the “Results” section: One for `setup`, one for `plot-dogs`, and one for `table-dogs`.
5. Write R code to load the packages and data in the `setup` chunk. Knit the document to see what it looks like.

<sup>4</sup>Knitting to PDF requires LaTeX, which you can install from tinytex. After you install LaTeX you can create truly beautifully typeset PDF documents.

<sup>5</sup>If you really want to do it, the best way is to create your template in the `.odt`, and then convert to `.docx`.

6. Write R code to create a simple `ggplot2` plot in the `plot-dogs` chunk related to your “Objective”. Knit the document to see what it looks like.
7. Write R code to create a simple summary table using `kable()` in the `table-dogs` chunk related to your “Objective”. Knit the document to see what it looks like.
8. Write an observation you made about the data from the plot and table in the “Conclusions” section. Include a **bold** text in this section.
9. Check out the Bootstrap themes and change your HTML theme to something else and add a Table of Contents.

## 5.8 Key Points

- A structured and standard project folder and file layout is the first step to having a reproducible data analysis project.
- Writing documents in R Markdown can reduce the time spent on manual tasks since results can be easily re-generated and inserted into the final document, improving reproducibility.
- `usethis` has several nifty helper functions for managing data analysis projects.
- Following a style guide and emphasizing readable code can lead to better quality code and to code that is more likely to be reproducible and reusable.
- Using Markdown to write documents is a great way to improve accessibility (since it is plain text only) and allows you to generate multiple types of output (HTML, PDF, slides, etc) from a single document source.

## 5.9 Additional learning resources and material

### For learning:

- Use other programming languages in an R Markdown document.
- Online book for R Markdown
- R Markdown chapter in the R for Data Science book.

### For help:

- RStudio helpful cheatsheets
- R Markdown cheatsheet (downloads a pdf file)
- R Markdown reference cheatsheet

Note: Source material for this chapter was modified from <https://rqawr-course.lwjohnst.com/>, as well as many other resources (see <https://rqawr-course.lwjohnst.com/license/>).





## Chapter 6

# Data Manipulation

### 6.1 Questions

- How can I read tabular data into a program?
- How can I select subsets of my data?
- How can I calculate new values?
- How can I tell what's gone wrong in my programs?
- How can I operate on subsets of my data?
- How can I save my results?

### 6.2 Motivation

TODO: where is data introduced?

TODO: clean version with no missing values and snake case column names, including month and year columns for some (all?) of the date columns.

TODO: I'm assuming that at some point this data will live in R package, so we can delay importing data to end of chapter.

The RStudio Viewer has an interface much like other spreadsheet programs you might have used. You can use this Viewer to look at the `dog_licenses` tibble with the `View()` function:

```
View(dog_licenses)
```

|   | row_number | animal_name | animal_gender | animal_birth_month | breed_name                         | borough |
|---|------------|-------------|---------------|--------------------|------------------------------------|---------|
| 1 | 533        | BONITA      | F             | 2013-05-01         | Unknown                            | Queens  |
| 2 | 548        | ROCKY       | M             | 2014-05-01         | Labrador Retriever Crossbreed      | Queens  |
| 3 | 622        | BULLY       | M             | 2010-07-01         | American Pit Bull Terrier/Pit Bull | Queens  |
| 4 | 633        | COCO        | M             | 2005-02-01         | Labrador Retriever                 | Queens  |
| 5 | 655        | SKI         | F             | 2012-09-01         | American Pit Bull Terrier/Pit Bull | Queens  |
| 6 | 872        | CHASE       | M             | 2013-11-01         | Shih Tzu                           | Queens  |
| 7 | 874        | CHEWY       | M             | 2014-09-01         | Shih Tzu                           | Queens  |
| 8 | 875        | CHASE       | M             | 2008-08-01         | Labrador Retriever                 | Queens  |

Showing 1 to 9 of 118,542 entries, 15 total columns

This viewer has some basic data manipulation features:

- **Arrange** You can change the order of the rows in the data based on the values in a column by clicking the up/down arrow next to the column name.
- **Filter** You can filter to include only rows which have a certain value in a column by first clicking the small funnel icon labelled “Filter”, then typing a desired value in the appropriate column.

Arrange and filter are known as data manipulation **verbs**. Individually, they describe a single simple manipulation of a dataset. It’s surprising how many questions you can answer using just these two basic verbs:

- How old is the oldest dog in this data? To answer we can arrange the `animal_birth_month` column in increasing order and see Jack, a Pug from Queens, was born in January 1999 (this license was issued in May 2015, making Jack at least 16 at the time). You’ll notice that there are other dogs with this same birthday.
- What range of license issue dates are in this data? Arrange `license_issued_date` once in increasing order and once in decreasing order, to find the issue dates range from 12th September 2015 to 31st December 2016.
- How many dogs licenses belong to dogs named Fido? Filter the `animal_name` column with `Fido`, and see “Showing ... of 12 entries” - so 12!
- and many more...

While these verbs are powerful in their own right, their real power comes from combining them. For example, we can answer the more complicated question “Which dog named Fido is the oldest?” by first filtering then arranging.

The Viewer in RStudio, however, has two huge limitations:

- It's a point and click interface. This means to repeat the same operation again you need to remember exactly the steps of point, clicking and typing you performed to get to your answer. Consequently, it's hard to share those steps unambiguously with someone else, and it's hard to save your results for future.
- The manipulation verbs in the viewer are limited. There is no way to rearrange the columns, add new variables or calculate summaries like counts or averages.

You'll start this chapter by overcoming this first limitation. You won't use the Viewer to arrange and filter, you'll learn to write code to do the same operations. Then you'll increase your vocabulary of data manipulation verbs to include:

- selecting variables,
- adding variables,
- summarizing rows, and
- performing these operations on subsets of the data.

Combining these verbs you'll be able to answer questions like:

- How long are licences issued for?
- What are the most popular breeds?
- What names are most popular for licensed dogs in New York? Does this vary geographically?
- When are dogs born?

TODO: update these question to reflect things that are actually done in this and later chapters.

To master data manipulation you need to master two pieces:

- How to describe the action you want with the data manipulation verbs individually. This is a language specific skill - in this chapter, you'll use the functions in the dplyr package.
- Identifying which verbs, and in which order to apply them, to answer a question of interest. This skill will translate across all technologies, but it takes a little longer to master.

### 6.2.1 Exercise: Point and click data manipulation

Using the RStudio Viewer answer the following questions:

- How many dog licenses belong to dogs named "Queen" that live in "Queens"?

## 6.3 Exploring data in the console

Let's take a look at the data in the console:

```
dog_licenses
```

```
## # A tibble: 118,542 x 15
##   row_number animal_name animal_gender animal_birth_month breed_name
##         <dbl> <chr>         <chr>         <dtm>         <chr>
## 1         533 BONITA         F         2013-05-01 00:00:00 Unknown
## 2         548 ROCKY         M         2014-05-01 00:00:00 Labrador ~
## 3         622 BULLY         M         2010-07-01 00:00:00 American ~
## 4         633 COCO         M         2005-02-01 00:00:00 Labrador ~
## 5         655 SKI          F         2012-09-01 00:00:00 American ~
## 6         872 CHASE         M         2013-11-01 00:00:00 Shih Tzu
## 7         874 CHEWY         M         2014-09-01 00:00:00 Shih Tzu
## 8         875 CHASE         M         2008-08-01 00:00:00 Labrador ~
## 9         893 MILEY         F         2008-07-01 00:00:00 Boxer
## 10        919 KENZI         F         2010-05-01 00:00:00 Schnauzer~
## # ... with 118,532 more rows, and 10 more variables: borough <chr>,
## #   zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## #   neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## #   congressional_district <dbl>, state_senatorial_district <dbl>,
## #   license_issued_date <date>, license_expired_date <date>
```

Notice that in contrast to the Viewer you only see the first 10 rows of the dataset, and just the first few columns. The number of columns you see depends on the width of your console, so you may see more or fewer than displayed here. You should also note some of the contents of the columns have been abbreviated. The ... at the end of some values in `breed_name` indicates these values have been truncated for display purposes.

You'll be using the `dplyr` package for data manipulation. Since it is part of the `tidyverse`, you'll need to load the `tidyverse` package to begin:

```
library(tidyverse)
```

### 6.3.1 Re-arranging rows

You can reorder rows of data with the `dplyr` function `arrange()`. The `arrange` function takes a tibble as its first argument and column names as the remaining arguments. The result will have the rows ordered in increasing value of the specified column. For example, to find the licenses belonging to the oldest dogs we arrange `dog_licenses` using the `animal_birth_month` column:

```
arrange(dog_licenses, animal_birth_month)
```

```
## # A tibble: 118,542 x 15
##   row_number animal_name animal_gender animal_birth_month breed_name
##   <dbl> <chr> <chr> <dtm> <chr>
## 1 15568 JACK M 1999-01-01 00:00:00 Pug
## 2 23695 KATTY F 1999-01-01 00:00:00 Chihuahua
## 3 101309 TOMMY M 1999-01-01 00:00:00 Unknown
## 4 1598 SARAH F 1999-01-01 00:00:00 West High~
## 5 8628 DOMINO M 1999-01-01 00:00:00 Labrador ~
## 6 15733 BRINKS M 1999-01-01 00:00:00 Yorkshire~
## 7 30419 LUCKY M 1999-01-01 00:00:00 Unknown
## 8 31348 MAGGIE F 1999-01-01 00:00:00 German Sh~
## 9 34685 COOKIE M 1999-01-01 00:00:00 Pomeranian
## 10 37194 DARCY F 1999-01-01 00:00:00 Cocker Sp~
## # ... with 118,532 more rows, and 10 more variables: borough <chr>,
## # zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## # neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## # congressional_district <dbl>, state_senatorial_district <dbl>,
## # license_issued_date <date>, license_expired_date <date>
```

You'll see Jack the Pug that lives in Queens, just like you did in the Viewer.

To arrange the rows by decreasing value, you need to wrap the column name in `desc()` (short for *descending* order). For instance to find the youngest dogs:

```
arrange(dog_licenses, desc(animal_birth_month))
```

```
## # A tibble: 118,542 x 15
##   row_number animal_name animal_gender animal_birth_month breed_name
##   <dbl> <chr> <chr> <dtm> <chr>
## 1 120352 MARLEY M 2016-12-01 00:00:00 Cocker Sp~
## 2 121981 MR. M 2016-12-01 00:00:00 Chihuahua~
## 3 120501 RORY M 2016-12-01 00:00:00 Unknown
## 4 122028 TAQUITO M 2016-12-01 00:00:00 Papillon
## 5 115820 REX M 2016-11-01 00:00:00 Maltese
## 6 121727 CHANDERBAL~ M 2016-11-01 00:00:00 Havanese
## 7 115777 ANGEL F 2016-11-01 00:00:00 Poodle, M~
## 8 118175 MASON M 2016-11-01 00:00:00 American ~
## 9 121601 TEDDY M 2016-11-01 00:00:00 Havanese
## 10 120995 LOLA F 2016-11-01 00:00:00 Morkie
## # ... with 118,532 more rows, and 10 more variables: borough <chr>,
## # zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## # neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## # congressional_district <dbl>, state_senatorial_district <dbl>,
## # license_issued_date <date>, license_expired_date <date>
```

As another example, to find the earliest issue date we can order by increasing `license_issued_date`:

```
arrange(dog_licenses, license_issued_date)
```

```
## # A tibble: 118,542 x 15
##   row_number animal_name animal_gender animal_birth_month breed_name
##         <dbl> <chr>         <chr>         <dtm>         <chr>
## 1           1 QUEEN           F           2013-04-01 00:00:00 Akita Cro~
## 2           2 CHEWBACCA       F           2012-06-01 00:00:00 Labrador ~
## 3           3 IAN             M           2006-01-01 00:00:00 Unknown
## 4           7 LOLA             F           2009-06-01 00:00:00 Maltese
## 5           4 PAIGE           F           2014-07-01 00:00:00 American ~
## 6           5 BUDDY           M           2008-04-01 00:00:00 Unknown
## 7           8 YOGI             M           2010-09-01 00:00:00 Boxer
## 8          10 MUNECA          F           2013-05-01 00:00:00 Beagle
## 9          27 BESS           F           2010-09-01 00:00:00 Beagle
## 10         26 BIGS           M           2004-12-01 00:00:00 American ~
## # ... with 118,532 more rows, and 10 more variables: borough <chr>,
## #   zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## #   neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## #   congressional_district <dbl>, state_senatorial_district <dbl>,
## #   license_issued_date <date>, license_expired_date <date>
```

The first row is the record with the earliest issue date, but we can't actually see that date because the column `license_issued_date` isn't being displayed due to space. One solution is to extract only the columns we are interested in, a manipulation known as selecting columns.

### 6.3.2 Exercise: Arranging character strings

Use `arrange()` to order the dog licenses by `animal_name` in increasing order. What does this tell you about the way R treats punctuation and numbers when dealing with alphabetical order?

## 6.4 How can I select subsets of my data?

Two verbs are used to subset data:

- `select()` to select columns
- `filter()` to select rows

You'll learn about these two functions in this section, along with learning about a way to chain together multiple operations on a dataset.

### 6.4.1 Selecting columns

The `select()` function in `dplyr` is used to extract some subset of columns (but keep all the rows) from a tibble. Just like `arrange()`, it takes a tibble as its first argument and column names as the remaining arguments. For example, to keep only the `animal_name` column:

```
select(dog_licenses, animal_name)
```

```
## # A tibble: 118,542 x 1
##   animal_name
##   <chr>
## 1 BONITA
## 2 ROCKY
## 3 BULLY
## 4 COCO
## 5 SKI
## 6 CHASE
## 7 CHEWY
## 8 CHASE
## 9 MILEY
## 10 KENZI
## # ... with 118,532 more rows
```

You can provide additional column names as arguments to keep additional specified columns, for example to keep `animal_name` and `breed_name`:

```
select(dog_licenses, animal_name, breed_name)
```

```
## # A tibble: 118,542 x 2
##   animal_name breed_name
##   <chr>         <chr>
## 1 BONITA      Unknown
## 2 ROCKY      Labrador Retriever Crossbreed
## 3 BULLY      American Pit Bull Terrier/Pit Bull
## 4 COCO      Labrador Retriever
## 5 SKI      American Pit Bull Terrier/Pit Bull
## 6 CHASE      Shih Tzu
## 7 CHEWY      Shih Tzu
## 8 CHASE      Labrador Retriever
## 9 MILEY      Boxer
## 10 KENZI     Schnauzer, Miniature
## # ... with 118,532 more rows
```

To return to finding the earliest issue date, you need to first arrange by increasing `license_issued_date` and then select the `license_issued_date` column. One approach is to store the result of the arrange step,

```
dog_by_date <- arrange(dog_licenses, license_issued_date)
```

Then apply the select step to this object:

```
select(dog_by_date, license_issued_date)
```

```
## # A tibble: 118,542 x 1
##   license_issued_date
##   <date>
## 1 2014-09-12
## 2 2014-09-12
## 3 2014-09-12
## 4 2014-09-12
## 5 2014-09-12
## 6 2014-09-12
## 7 2014-09-12
## 8 2014-09-13
## 9 2014-09-13
## 10 2014-09-13
## # ... with 118,532 more rows
```

There are lots of shortcuts you can use with `select()` to avoid having to type out all the variables you want to keep. For example, you can ask for all the columns that start with a certain string:

```
select(dog_licenses, starts_with("Animal"))
```

```
## # A tibble: 118,542 x 3
##   animal_name animal_gender animal_birth_month
##   <chr>      <chr>      <dtm>
## 1 BONITA    F          2013-05-01 00:00:00
## 2 ROCKY     M          2014-05-01 00:00:00
## 3 BULLY     M          2010-07-01 00:00:00
## 4 COCO      M          2005-02-01 00:00:00
## 5 SKI       F          2012-09-01 00:00:00
## 6 CHASE     M          2013-11-01 00:00:00
## 7 CHEWY     M          2014-09-01 00:00:00
## 8 CHASE     M          2008-08-01 00:00:00
## 9 MILEY     F          2008-07-01 00:00:00
## 10 KENZI    F          2010-05-01 00:00:00
## # ... with 118,532 more rows
```

Take a look in the “Useful functions” section of the `select()` help page for a complete list:

```
?dplyr::select
```



### 6.4.2 Exercise: Find the latest license issue date

Combine `arrange()` and `select()` to confirm the last issue date in this dataset is 31st December 2016.

### 6.4.3 Combining operations with the pipe `%>%`

You've seen you can combine data manipulation steps to do more complicated tasks, but so far you've done so by saving an intermediate object, in our previous example the object `dog_by_date`:

```
dog_by_date <- arrange(dog_licenses, license_issued_date)
select(dog_by_date, license_issued_date)
```

The pipe, `%>%`, is an operator that allows you to chain together operations without intermediate objects and maintain readability. The name, “pipe”, comes from the plumbing kind of pipe, not the smoking kind, and references the idea of objects flowing out of one function and into another. Let's just look at the first step in our manipulation:

```
arrange(dog_licenses, license_issued_date)
```

With the pipe this can be rewritten as:

```
dog_licenses %>% arrange(license_issued_date)
```

The pipe takes the object on the left hand side and passes it as the first argument to the function on the right hand side. So, here the `dog_licenses` dataset is passed to the first argument of `arrange()`. Inside `arrange()` we can then list any additional arguments as we normally would.

The pipe works very nicely with the data manipulation verbs because every verb expects a tibble as its first argument and returns a tibble. This means the result of one operation is easily piped into the next operation, allowing you to chain together multiple steps. For instance, piping the result of the `arrange` step above into the `select()` function:

```
dog_licenses %>%
  arrange(license_issued_date) %>%
  select(license_issued_date)
```

When you see the pipe, read it as “and then”. So, the above code would be read:

Take the `dog_licenses` data, **and then**  
arrange the rows by the `license_issued_date`, **and then**  
select the column `license_issued_date`.

The result is code that matches very closely how we might describe the steps we performed in natural language. It's so natural that for the remainder of the

chapter we'll use the pipe when combining data manipulation steps.

#### 6.4.4 Exercise: Reading aloud

Read the following code aloud to your neighbor (or cat, dog, or rubber duck). Remember to pronounce `%>%` as “and then”.

```
dog_licenses %>%
  arrange(license_issued_date) %>%
  select(license_expired_date)
```

What question might it answer?

#### 6.4.5 Exercise: Using the pipe

Re-write this snippet of code to use the pipe:

```
select(dog_licenses, animal_name, breed_name)
```

Use the pipe to re-write this snippet of code to avoid the intermediate variable:

```
name_and_breed <- select(dog_licenses, animal_name, breed_name)
arrange(name_and_breed, breed_name)
```

#### 6.4.6 Filtering to keep a subset of rows

The function to filter rows of a tibble is `filter()`. Like `arrange()` and `select()`, its first argument is a tibble. The remaining arguments describe which rows to **keep**. The rows to keep are specified with a logical expression - something that is either `TRUE` or `FALSE`. The rows where this expression is `TRUE` will be returned.

One of the simplest kinds of logical expression is a test for equality with the `==` operator. For example, to keep the rows where `animal_name` is `BRUNO` you could do:

```
dog_licenses %>% filter(animal_name == "BRUNO")
```

```
## # A tibble: 272 x 15
##   row_number animal_name animal_gender animal_birth_month breed_name
##   <dbl> <chr> <chr> <dtm> <chr>
## 1 12001 BRUNO M 2010-05-01 00:00:00 American ~
## 2 27228 BRUNO M 2013-07-01 00:00:00 Jack Russ~
## 3 68192 BRUNO M 2002-01-01 00:00:00 Shih Tzu
## 4 70175 BRUNO M 2015-12-01 00:00:00 Chihuahua~
## 5 120562 BRUNO M 2016-05-01 00:00:00 Labrador ~
```

```
## 6      3915 BRUNO      M      2014-03-01 00:00:00 Doberman ~
## 7      9614 BRUNO      M      2014-12-01 00:00:00 Boxer
## 8     15606 BRUNO      M      2015-02-01 00:00:00 French Bu~
## 9     32742 BRUNO      M      2014-03-01 00:00:00 Maltipoo
## 10    34299 BRUNO      M      2003-01-01 00:00:00 Unknown
## # ... with 262 more rows, and 10 more variables: borough <chr>,
## #   zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## #   neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## #   congressional_district <dbl>, state_senatorial_district <dbl>,
## #   license_issued_date <date>, license_expired_date <date>
```

You could read this code aloud as:

Take the `dog_licenses` data, **and then**  
 filter for only rows when `animal_name` is equal to "BRUNO"

The logical expression, `animal_name == "BRUNO"`, will be `TRUE` when the value of the `animal_name` column is exactly equal to the character string "BRUNO" - any differences in characters, case, or whitespace will result in `FALSE`.

Above, each value within a column was compared against the same fixed string ("BRUNO"). You can also compare the values from two columns against each other by including a column name on each side of the `==` sign. In programming terms this is known as element-wise comparison. For example, `license_issued_date == animal_birth_month` will return `TRUE` for a row only if for that row the date the license was issued is the exact same date as the birth month for the dog. If you take a look:

```
dog_licenses %>%
  filter(license_issued_date == animal_birth_month)

## # A tibble: 0 x 15
## # ... with 15 variables: row_number <dbl>, animal_name <chr>,
## #   animal_gender <chr>, animal_birth_month <dtm>, breed_name <chr>,
## #   borough <chr>, zip_code <dbl>, community_district <dbl>,
## #   census_tract_2010 <dbl>, neighborhood_tabulation_area <chr>,
## #   city_council_district <dbl>, congressional_district <dbl>,
## #   state_senatorial_district <dbl>, license_issued_date <date>,
## #   license_expired_date <date>
```

There is no output (apart from the column names), which means that no rows satisfy this criteria.

Remember that a string is surrounded by quotes while a column name is not. When you are reading code, look for the quotes to figure out if the comparison is to a string, or (by the absence of quotes) to the strings within a column. Use the same strategy to figure out where the quotes should be in your own code, but beware: misplacing quotes often won't result in an error, but instead a result that you weren't expecting. For example, I might be interested in the licenses

issued to male dogs and try:

```
dog_licenses %>%
  filter("animal_gender" == "M")
```

```
## # A tibble: 0 x 15
## #   ... with 15 variables: row_number <dbl>, animal_name <chr>,
## #     animal_gender <chr>, animal_birth_month <dtm>, breed_name <chr>,
## #     borough <chr>, zip_code <dbl>, community_district <dbl>,
## #     census_tract_2010 <dbl>, neighborhood_tabulation_area <chr>,
## #     city_council_district <dbl>, congressional_district <dbl>,
## #     state_senatorial_district <dbl>, license_issued_date <date>,
## #     license_expired_date <date>
```

The result has zero rows, which would suggest there are no such licenses, but in fact this is the answer to a different question. Can you see what is wrong with the code? By surrounding `animal_gender` in quotes, R has interpreted the comparison as: is the string "animal\_gender" equal to the string "M"? The answer is FALSE, and no rows are returned. I actually wanted to compare the `animal_gender` column to the string "M", so `animal_gender` should have no quotes around it:

```
dog_licenses %>%
  filter(animal_gender == "M")
```

```
## # A tibble: 64,770 x 15
##   row_number animal_name animal_gender animal_birth_month breed_name
##   <dbl> <chr> <chr> <dtm> <chr>
## 1     548 ROCKY      M      2014-05-01 00:00:00 Labrador ~
## 2     622 BULLY      M      2010-07-01 00:00:00 American ~
## 3     633 COCO       M      2005-02-01 00:00:00 Labrador ~
## 4     872 CHASE      M      2013-11-01 00:00:00 Shih Tzu
## 5     874 CHEWY      M      2014-09-01 00:00:00 Shih Tzu
## 6     875 CHASE      M      2008-08-01 00:00:00 Labrador ~
## 7     976 APOLLO     M      2014-10-01 00:00:00 American ~
## 8    1297 JERRY      M      2009-06-01 00:00:00 Labrador ~
## 9    2133 SIMON      M      2010-12-01 00:00:00 Havanese
## 10   2289 BUDDY      M      2012-06-01 00:00:00 Labrador ~
## #   ... with 64,760 more rows, and 10 more variables: borough <chr>,
## #     zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## #     neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## #     congressional_district <dbl>, state_senatorial_district <dbl>,
## #     license_issued_date <date>, license_expired_date <date>
```

The operator, `==` (you can read as “is equal to”, or simply “equals”), is a specific kind of comparison. Other comparisons include:

| Operator | Meaning                  |
|----------|--------------------------|
| <        | less than                |
| >        | greater than             |
| <=       | less than or equal to    |
| >=       | greater than or equal to |
| !=       | not equal to             |

### 6.4.7 Exercise: Enterprising dogs

Are there any dogs called “Spock”, “Picard”, or “Janeway”?

These are Star Trek characters, can you find any dogs with a name from one of your favorite books, TV shows, or movies?

### 6.4.8 Exercise: Expired Licenses

This code creates a variable that contains the date for the start of the year 2016:

```
start_of_2016 <- as.Date("2016-01-01")
```

Use `filter()` with this variable to find:

- The dog licenses that were issued before 2016
- The dog licenses that expire before 2016

### 6.4.9 More complicated expressions

Logical expressions can be combined with logical operators to construct more complicated expressions. For example, the AND operator, `&`, returns `TRUE` only if the expressions on both sides of it are `TRUE`. For example, you’ve seen, you can find the licenses to dogs called “BRUNO”:

```
dog_licenses %>%
  filter(animal_name == "BRUNO")
```

```
## # A tibble: 272 x 15
```

```
##   row_number animal_name animal_gender animal_birth_month breed_name
##   <dbl> <chr> <chr> <dtm> <chr>
## 1 12001 BRUNO M 2010-05-01 00:00:00 American ~
## 2 27228 BRUNO M 2013-07-01 00:00:00 Jack Russ~
## 3 68192 BRUNO M 2002-01-01 00:00:00 Shih Tzu
## 4 70175 BRUNO M 2015-12-01 00:00:00 Chihuahua~
## 5 120562 BRUNO M 2016-05-01 00:00:00 Labrador ~
## 6 3915 BRUNO M 2014-03-01 00:00:00 Doberman ~
```

```
## 7      9614 BRUNO      M      2014-12-01 00:00:00 Boxer
## 8     15606 BRUNO      M      2015-02-01 00:00:00 French Bu~
## 9     32742 BRUNO      M      2014-03-01 00:00:00 Maltipoo
## 10    34299 BRUNO      M      2003-01-01 00:00:00 Unknown
## # ... with 262 more rows, and 10 more variables: borough <chr>,
## #   zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## #   neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## #   congressional_district <dbl>, state_senatorial_district <dbl>,
## #   license_issued_date <date>, license_expired_date <date>
```

And you could find the dog licenses issued to dogs that live in Brooklyn:

```
dog_licenses %>%
  filter(borough == "Brooklyn")
```

```
## # A tibble: 29,334 x 15
##   row_number animal_name animal_gender animal_birth_month breed_name
##   <dbl> <chr> <chr> <dtm> <chr>
## 1      2895 FUDGE      M      2014-07-01 00:00:00 American ~
## 2       4057 STAR      F      2011-01-01 00:00:00 Poodle
## 3     74463 MUNECA      F      2011-09-01 00:00:00 Chihuahua~
## 4     76232 KATTY      F      2002-06-01 00:00:00 Chihuahua
## 5     85113 SHADOW      M      2015-03-01 00:00:00 American ~
## 6     85997 SPARKIE      M      2013-08-01 00:00:00 Maltese C~
## 7     92451 SNOW      M      2014-07-01 00:00:00 Maltese
## 8    104256 BELLA      F      2016-07-01 00:00:00 Maltese
## 9     10389 SPARKLE      F      2006-01-01 00:00:00 Schnauzer~
## 10    82492 UNKNOWN      F      2015-11-01 00:00:00 Pomeranian
## # ... with 29,324 more rows, and 10 more variables: borough <chr>,
## #   zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## #   neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## #   congressional_district <dbl>, state_senatorial_district <dbl>,
## #   license_issued_date <date>, license_expired_date <date>
```

If you want to find the licenses that are to dogs named Bruno in Brooklyn, you could combine the two logical statements with `&`:

```
dog_licenses %>%
  filter((animal_name == "BRUNO") & (borough == "Brooklyn"))
```

```
## # A tibble: 55 x 15
##   row_number animal_name animal_gender animal_birth_month breed_name
##   <dbl> <chr> <chr> <dtm> <chr>
## 1      8496 BRUNO      M      2005-04-01 00:00:00 American ~
## 2     65466 BRUNO      M      2012-11-01 00:00:00 Golden Re~
## 3    115999 BRUNO      M      2009-01-01 00:00:00 Pug
## 4    118963 BRUNO      M      2005-01-01 00:00:00 Unknown
## 5     10505 BRUNO      M      2006-04-01 00:00:00 Bull Dog,~
```

```
## 6      14444 BRUNO      M      2007-01-01 00:00:00 Cocker Sp~
## 7      14690 BRUNO      M      2014-04-01 00:00:00 Shih Tzu
## 8      47454 BRUNO      M      2011-06-01 00:00:00 Pug
## 9      58918 BRUNO      M      2006-04-01 00:00:00 Bull Dog,~
## 10     105964 BRUNO      M      2016-04-01 00:00:00 Bull Dog,~
## # ... with 45 more rows, and 10 more variables: borough <chr>,
## #   zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## #   neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## #   congressional_district <dbl>, state_senatorial_district <dbl>,
## #   license_issued_date <date>, license_expired_date <date>
```

The parentheses around each logical expression are optional, but can help visually to group the components to `&`, especially if those logical expressions get more complicated.

If you want to create an “or” type expression, like “licenses to dogs named BRUNO **or** dogs named BRUCE”, you need to combine two comparisons with the OR operator, `|`.

```
dog_licenses %>%
  filter((animal_name == "BRUNO") | (animal_name == "BRUCE"))
```

If you find yourself combining lots of comparisons on the same column with `|`, like dogs named BRUNO, BRUCE or BRADY:

```
dog_licenses %>%
  filter((animal_name == "BRUNO") | (animal_name == "BRUCE") | (animal_name == "BRADY"))
```

You can save a lot of typing with `%in%`:

```
dog_licenses %>%
  filter(animal_name %in% c("BRUNO", "BRUCE", "BRADY"))
```

On the right hand side of `%in%` the function `c()`, combines many single values into a vector. `%in%` will return `TRUE` for an element of the left hand side if it is contained in the vector on the right hand side.

#### 6.4.10 Exercise: Expired Licenses

This code creates two variables that contain the dates for the start and end of the year 2016:

```
start_of_2016 <- as.Date("2016-01-01")
endt_of_2016 <- as.Date("2016-12-31")
```

Use `filter()` with these variables to find the dog licenses that expire during 2016.

## 6.5 How can I calculate new values?

So far, you've been manipulating the columns that already exist in a dataset, but what if you want to add new ones? The function `mutate()` handles this kind of operation.

To see how this works let's start with a logical expression: `animal_name == "CHASE"`. If you used this with `filter()`, you would get all the rows back where the license was issued to a dog named CHASE. Let's say instead of subsetting the data, you want to add a column called `called_chase` that contained the TRUE and FALSE result. You might do this for example, if you are interested in comparing the two groups of dogs, rather than just keeping one of them. With `mutate()` after passing in the data, you pass named arguments, where the name is the name you desire for the new column, and its value is the way to calculate it:

```
dog_licenses %>%
  mutate(called_chase = animal_name == "CHASE") %>%
  select(animal_name, called_chase)
```

```
## # A tibble: 118,542 x 2
##   animal_name called_chase
##   <chr>         <lgl>
## 1 BONITA      FALSE
## 2 ROCKY       FALSE
## 3 BULLY       FALSE
## 4 COCO        FALSE
## 5 SKI         FALSE
## 6 CHASE       TRUE
## 7 CHEWY       FALSE
## 8 CHASE       TRUE
## 9 MILEY       FALSE
## 10 KENZI      FALSE
## # ... with 118,532 more rows
```

Take `dog_licences`, **and then**,  
`mutate` to add a column called `called_chase` which is the result of  
 testing whether `animal_name` is exactly "CHASE", **and then**,  
`select` the columns `animal_name` and `called_chase`.

The `select` statement isn't crucial to the calculation here, but it does help me draw your attention to the columns that were involved in this step.

This can be a useful intermediate step in filtering, since it gives you a chance to examine the logical statement before using it to filter:

```
dog_licenses %>%
  mutate(called_chase = animal_name == "CHASE") %>%
```



```
filter(called_chase)
```

```
## # A tibble: 126 x 16
##   row_number animal_name animal_gender animal_birth_month breed_name
##   <dbl> <chr>      <chr>      <dtm>      <chr>
## 1      872 CHASE      M      2013-11-01 00:00:00 Shih Tzu
## 2      875 CHASE      M      2008-08-01 00:00:00 Labrador ~
## 3     32652 CHASE      M      2013-09-01 00:00:00 Yorkshire~
## 4     42125 CHASE      M      2015-08-01 00:00:00 Chihuahua
## 5     109847 CHASE      M      2013-04-01 00:00:00 Terrier m~
## 6     114557 CHASE      M      2009-07-01 00:00:00 Siberian ~
## 7      33434 CHASE      M      2014-08-01 00:00:00 Schnauzer~
## 8      45142 CHASE      M      2005-01-01 00:00:00 Unknown
## 9       2528 CHASE      M      2011-10-01 00:00:00 Lhasa Apso
## 10      3426 CHASE      M      2014-03-01 00:00:00 Yorkshire~
## # ... with 116 more rows, and 11 more variables: borough <chr>,
## #   zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## #   neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## #   congressional_district <dbl>, state_senatorial_district <dbl>,
## #   license_issued_date <date>, license_expired_date <date>,
## #   called_chase <lgl>
```

Take a closer look at the argument to mutate:

```
called_chase = animal_name == "CHASE"
```

On the left of the `=` is the argument name, `called_chase`. This is a name you choose - it should be descriptive and follow good style. On the right of the `=` is an expression that must either return as many values as there are rows, or a single value. Here, the logical expression involves the column `animal_name` so it returns as many values as there are rows.

If you would prefer the values to be something other than `TRUE` or `FALSE`, instead maybe you want them to be something like "called chase" or "not called chase", we would need to use the `ifelse()`.

A call to `ifelse()` takes the form:

```
ifelse(test, yes, no)
```

Where `test` is a logical expression, `yes` the value for the elements that return `TRUE`, and `no` the value for the elements that return `FALSE`. (Both `yes` and `no` could be other column names, in which case, the corresponding element of `yes` would be returned for `TRUE` elements).

```
dog_licenses %>%
  mutate(
    called_chase = ifelse(animal_name == "CHASE", "called chase", "not called chase")
  ) %>%
```

```
select(animal_name, called_chase)
```

```
## # A tibble: 118,542 x 2
##   animal_name called_chase
##   <chr>         <chr>
## 1 BONITA      not called chase
## 2 ROCKY       not called chase
## 3 BULLY       not called chase
## 4 COCO        not called chase
## 5 SKI         not called chase
## 6 CHASE       called chase
## 7 CHEWY       not called chase
## 8 CHASE       called chase
## 9 MILEY       not called chase
## 10 KENZI      not called chase
## # ... with 118,532 more rows
```

FIXME: add diagram showing how ifelse works

You can perform multiple mutate steps at once by passing more arguments to `mutate()`, so an alternative way of writing the above code (with more keystrokes, but with lines that are shorter) would be:

```
dog_licenses %>%
  mutate(
    is_chase = animal_name == "CHASE",
    called_chase = ifelse(is_chase, "called chase", "not called chase")
  ) %>%
  select(animal_name, is_chase, called_chase)
```

```
## # A tibble: 118,542 x 3
##   animal_name is_chase called_chase
##   <chr>       <lgl>    <chr>
## 1 BONITA     FALSE    not called chase
## 2 ROCKY      FALSE    not called chase
## 3 BULLY      FALSE    not called chase
## 4 COCO       FALSE    not called chase
## 5 SKI        FALSE    not called chase
## 6 CHASE      TRUE     called chase
## 7 CHEWY      FALSE    not called chase
## 8 CHASE      TRUE     called chase
## 9 MILEY      FALSE    not called chase
## 10 KENZI     FALSE    not called chase
## # ... with 118,532 more rows
```

Notice that the computation for the `called_chase` column refers to the `is_chase` column. The arguments to `mutate` are computed in order, so columns created later in the same `mutate()` can refer to columns created

earlier.

Arithmetic is another common operation that returns as many elements as there are rows. For instance we could see how long licenses are issued for:

```
dog_licenses %>%
  mutate(license_duration = license_expired_date - license_issued_date) %>%
  select(license_duration)
```

```
## # A tibble: 118,542 x 1
##   license_duration
##   <drtn>
## 1 1118 days
## 2 1826 days
## 3  697 days
## 4 1096 days
## 5 1826 days
## 6  731 days
## 7  731 days
## 8 1097 days
## 9  421 days
## 10 402 days
## # ... with 118,532 more rows
```

The output shows us licenses aren't issued for a standard time period. In these first rows, there are some licenses issued for a whole number of years: 2 (731 days), 3 (1097 days) and 5 (1826 days). However, others seem to be for fractions of years like 421 days.

TODO: provide link to RStudio data mini cheatsheet with other functions that are useful with mutate.

You can use `mutate()` with operations that give one number based on all the rows:

```
dog_licenses %>%
  mutate(
    license_duration = license_expired_date - license_issued_date,
    avg_duration = mean(license_duration)) %>%
  select(license_duration, avg_duration)
```

```
## # A tibble: 118,542 x 2
##   license_duration avg_duration
##   <drtn>          <drtn>
## 1 1118 days      467.321 days
## 2 1826 days      467.321 days
## 3  697 days      467.321 days
## 4 1096 days      467.321 days
## 5 1826 days      467.321 days
```

```
## 6 731 days      467.321 days
## 7 731 days      467.321 days
## 8 1097 days     467.321 days
## 9 421 days      467.321 days
## 10 402 days     467.321 days
## # ... with 118,532 more rows
```

You'll get back the original number of rows, but the single value will be repeated in all of them - on average licenses are issued for 467.321 days. You'll see a different verb, `summarise()` that collapses many rows into one later in this chapter.

### 6.5.1 Exercise: Ages of dogs

Use `mutate()` to add a column `age_at_issue` that contains the dogs approximate age on the day the license was issued.

### 6.5.2 Exercise: Unknown breeds

Use `mutate()` along with `ifelse()` to create a column `breed` that takes values "unknown" if `breed_name` is "unknown" and "known" otherwise.

*Extra challenge* Can you figure out if the licenses issued to unknown breed dogs are of longer or shorter duration on average than known breed dogs?

### 6.5.3 Exercise: Name length

The function `str_length()` in the `stringr` package finds the length of character strings. Replace the `___` in following code to add a column called `name_length` that contains the length of the dog's name:

```
dog_licenses %>%
  ___(___ = stringr::str_length(animal_name))
```

Now add an `arrange()` step to find the licenses issued to dogs with the longest names?

## 6.6 How can I tell what's gone wrong in my programs?

TODO: The errors shown in the markdown are way more informative than those in the Console. Try to get the error displaying in the book like they do for someone in the console?

## 6.6. HOW CAN I TELL WHAT'S GONE WRONG IN MY PROGRAMS?101

Let me share an interaction my (CVW's) husband had at a Trader Joes (a small specialized supermarket) soon after we arrived in the USA from New Zealand.

Josh: "Do you sell *bat-trees*?"

Store-person: "What?"

Josh: "Do you sell *bat-trees*?"

Store person: "Huh? *Bat...trees*?"

Josh: "Do you sell *bat-er-ries*?"

Store-person: "Oh...you mean batteries! No. We don't sell batteries."

This is a pretty accurate analogy for what it feels like when you are learning R. You know what you want, but you have to ask for it in a way R understands. When R doesn't understand you, or when R can't give you what you want, you'll get an error. You'll know when you get one in R because the only output you see will start with **Error**.

```
buy_batteries(store = "Trader Joes")
```

```
## Error in buy_batteries(store = "Trader Joes"): could not find function "buy_batteries"
```

It's also a good illustration of the two kinds of problems that occur: syntax errors and runtime errors. Syntax errors are like the "What? Huh?" moments. R doesn't understand what you are asking it to do, because something about the way you are asking doesn't conform to what R expects and it will not even try to run your code. Runtime errors are more like the "No, I can't help you" moments. R understands what you are asking and runs your code, but during the run something goes wrong and R has to stop running the code.

You generally want to make sure you've ruled out syntax errors before assuming it's a runtime error. Unfortunately, R doesn't distinguish these two types of error in its output, so we'll discuss some of the most common examples in the following sections.

This analogy is also a reminder that sometimes you'll have to repeat yourself. With R, giving the exact same instruction should always result in the exact same error, but it's not uncommon, even for longtime R users, to run and edit a line of code multiple times before it runs without error.

There is one flaw in this analogy – R isn't a real person. So, if you need to vent *your* frustration by cursing it, insulting its parentage or storming off, it's fine, no one's feelings will be hurt. Of course, it won't change R's response...

### 6.6.1 Common syntax errors

Syntax errors occur when your code can't be broken into its component pieces by R. For example, R expects the arguments to a function to start after the

opening parenthesis (`(`), be separated by a comma (`,`) and finish at the closing parenthesis (`)`). When R sees this code:

```
filter(dog_licenses breed_name == "Finnish Lapphund")
```

```
## Error: <text>:1:21: unexpected symbol
## 1: filter(dog_licenses breed_name
##                               ^
```

The missing comma means R can't figure out where the first argument to `filter()` ends: is it after `dog_licenses`, after `breed_name`, or after `==`?

Take a closer look at the error message. Error messages always begin with **Error**, then optionally the name of the function that returned an error (not in this example), followed by a `:`, and some description of the error that occurred. The message **unexpected symbol** is one common kind of syntax error - in this case R encountered some code where it was expecting a comma or closing parenthesis. We can fix it by putting in the missing comma:

```
filter(dog_licenses, breed_name == "Finnish Lapphund")
```

```
## # A tibble: 1 x 15
##   row_number animal_name animal_gender animal_birth_month breed_name
##   <dbl> <chr>      <chr>      <dtm>      <chr>
## 1    118408 FREDDIE      M          2011-11-01 00:00:00 Finnish L~
## # ... with 10 more variables: borough <chr>, zip_code <dbl>,
## #   community_district <dbl>, census_tract_2010 <dbl>,
## #   neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## #   congressional_district <dbl>, state_senatorial_district <dbl>,
## #   license_issued_date <date>, license_expired_date <date>
```

Syntax errors are usually the result of typos. Some things to keep an eye out for:

- If you are modelling your code on an example, pay very close attention to the punctuation: commas `,`, parenthesis `(, )`, brackets `[, ]`, and quotes `"`, `'`. Every opening parenthesis, bracket, or quote needs a matching closing one, and they must be closed in the reverse order they were opened.
- New lines don't matter if they happen between arguments, or after pipe operators, but they can be problematic in other locations.
- R ignores other whitespace (spaces or tabs) unless it's inside a character string (i.e. inside quotes), so different spacing shouldn't be the cause of an error, but it is a good idea to follow good style for spacing. TODO: link to style guide section.

### 6.6.2 Exercise: Syntax errors

Fix these **syntax** errors.

- `dog_licenses %>%  
 arrange(desc(license_expired_date)))`

```
## Error: <text>:2:38: unexpected ' '
## 1: dog_licenses %>%
## 2:   arrange(desc(license_expired_date)))
##                                     ^
```

- `dog_licenses %>%  
 mutate(  
 month_born = lubridate::month(animal_birth_month)  
 year_born = lubridate::year(animal_birth_month))`

```
## Error: <text>:4:5: unexpected symbol
## 3:   month_born = lubridate::month(animal_birth_month)
## 4:   year_born
##      ^
```

- `dog_licenses %>%  
 filter(animal_name == "BRUNO)`

```
## Error: <text>:2:25: unexpected INCOMPLETE_STRING
## 1: dog_licenses %>%
## 2:   filter(animal_name == "BRUNO)
##                                     ^
```

(If you run this code in the Console, you might not get an error, but you should see a + on a new line, a signal that R is waiting for more input and a clue that there is something missing in this code).

- `dog_licenses  
 %>% filter(animal_gender == "M")`

```
## Error: <text>:2:3: unexpected SPECIAL
## 1: dog_licenses
## 2:   %>%
##      ^
```

### 6.6.3 Common runtime errors

Runtime errors come in an infinite number of flavors because there are so many ways that you ask for that might be impossible to do. For example you might

try to do arithmetic with character strings:

```
"apple" + "banana"
```

```
## Error in "apple" + "banana": non-numeric argument to binary operator
```

Or try to filter with something that isn't a logical:

```
dog_licenses %>% filter(animal_name)
```

```
## Error: Argument 2 filter condition does not evaluate to a logical vector
```

Perhaps the most common runtime error is of the form `Error: object not found`:

```
an_object_i_dont_have
```

```
## Error in eval(expr, envir, enclos): object 'an_object_i_dont_have' not found
```

This is R complaining that you've asked it to operate on an object that it doesn't know about. Often this is actually a typo in disguise, for example you've misspelled the name of the object,

```
my_object <- 12
my_objet
```

```
## Error in eval(expr, envir, enclos): object 'my_objet' not found
```

you've used the wrong case,

```
My_object
```

```
## Error in eval(expr, envir, enclos): object 'My_object' not found
```

or you've forgotten that you've used a separator

```
myobject
```

```
## Error in eval(expr, envir, enclos): object 'myobject' not found
```

This error also often arises when you forgot quotes around strings. For example, if we want all the dogs that are male, and try

```
dog_licenses %>% filter(animal_gender == M)
```

```
## Error: object 'M' not found
```

you get an error because R is looking for an object called `M` to compare to the values in the column called `animal_gender`. What we really wanted to do was compare the values in `animal_gender` to the string `"M"`:

```
dog_licenses %>% filter(animal_gender == "M")
```

```
## # A tibble: 64,770 x 15
```

```
##   row_number animal_name animal_gender animal_birth_month breed_name
```



## 6.6. HOW CAN I TELL WHAT'S GONE WRONG IN MY PROGRAMS?105

```
##           <dbl> <chr>           <chr>           <dtm>           <chr>
## 1           548 ROCKY             M             2014-05-01 00:00:00 Labrador ~
## 2           622 BULLY             M             2010-07-01 00:00:00 American ~
## 3           633 COCO              M             2005-02-01 00:00:00 Labrador ~
## 4           872 CHASE             M             2013-11-01 00:00:00 Shih Tzu
## 5           874 CHEWY             M             2014-09-01 00:00:00 Shih Tzu
## 6           875 CHASE             M             2008-08-01 00:00:00 Labrador ~
## 7           976 APOLLO            M             2014-10-01 00:00:00 American ~
## 8          1297 JERRY             M             2009-06-01 00:00:00 Labrador ~
## 9          2133 SIMON             M             2010-12-01 00:00:00 Havanese
## 10         2289 BUDDY             M             2012-06-01 00:00:00 Labrador ~
## # ... with 64,760 more rows, and 10 more variables: borough <chr>,
## #   zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## #   neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## #   congressional_district <dbl>, state_senatorial_district <dbl>,
## #   license_issued_date <date>, license_expired_date <date>
```

### 6.6.4 Exercise: object not found

Fix these `object not found` errors. (Hint: the names of the objects or variables being created should give you a clue to the intent of the code)

```
• dog_licenses %>%
  mutate(year_issued = lubridate::year(Liscenceissuedate))

## Error in lubridate::year(Liscenceissuedate): object 'Liscenceissuedate' not found

• dogs_named_bruno <- dog_licenses %>%
  filter(animal_name == BRUNO)

## Error: object 'BRUNO' not found
```

### 6.6.5 Warnings and messages

There are two other kinds of alerts R can give: warnings and messages. These can both appear in the console with the same color as an error, but they are informational as opposed to fatal.

Messages are purely informational, for example when you read data in with `read_csv()` you get a message that describes the columns and their data types as parsed by the function:

```
sites <- read_csv("site.csv")

## Parsed with column specification:
## cols(
```

```
##   site_id = col_character(),
##   latitude = col_double(),
##   longitude = col_double()
## )
```

Warnings generally alert you that something was slightly unexpected but that R recovered and gave you a result anyway. Poorly formatted CSV files will often result in warnings from `read_csv()`:

```
bad_csv <- "
id,
1, 2
2, 1
3, 5
"
bad <- read_csv(bad_csv)
```

```
## Warning: Missing column names filled in: 'X2' [2]
```

Here there was a missing column name. `read_csv()` still returns an object but the warning alerts you that it made some assumption to get it, i.e. that it made up a column name:

```
bad

## # A tibble: 3 x 2
##       id      X2
##   <dbl> <dbl>
## 1     1     2
## 2     2     1
## 3     3     5
```

Warnings don't **stop** you from proceeding, but they should alert you to question whether you **should be** proceeding.

### 6.6.6 What do I do when I get an Error I can't fix?

- Check that the error is reproducible. Restart R with a clean slate and re-run your code up to and including the code that gives the error. This is the R version of the classic tech advice to “turn it off, then turn it on again”. TODO: link to reproducibility chapter.
- Try searching for it online: for example, searching for “R unexpected string constant” lead me to the question “Error: unexpected symbol/input/string constant/numeric constant/SPECIAL in my code” on StackOverflow which gives some great examples of ways this error might arise.

- Ask for help. You are most likely to get help when you can provide a reproducible example. Stack Overflow has detailed instruction on how to create a minimal reproducible example when asking a question to increase the chances that the question receives a specific and helpful answer. The key principles listed on the website recommends that an answer follows these guidelines:
  - Minimal – Use as little code as possible that still produces the same problem
  - Complete – Provide all parts someone else needs to reproduce your problem in the question itself
  - Reproducible – Test the code you’re about to provide to make sure it reproduces the problem

## 6.7 How can I operate on subsets of my data?

The syntax for `summarise()` is the same as `mutate()` but it expects operations that reduce all rows down to one row. Recall from `mutate()` that this code added the average license duration to every row of the data:

```
dog_licenses %>%
  mutate(license_duration = license_expired_date - license_issued_date) %>%
  mutate(avg_duration = mean(license_duration))
```

You actually saw this in one `mutate()` statement, but I’ve separated out the line that calculates the average so it’s easier to see the difference with `summarise()`. See what happens when you switch out the final `mutate()` with `summarise()`:

```
dog_licenses %>%
  mutate(license_duration = license_expired_date - license_issued_date) %>%
  summarise(avg_duration = mean(license_duration))
```

```
## # A tibble: 1 x 1
##   avg_duration
##   <drtn>
## 1 467.321 days
```

Instead of the the one value repeated on every row, we get a new tibble with only one row, and a single column that corresponds to our requested summary.

Any function that takes many values and reduces them to one is a good candidate for `summarise()`, for example we could find the shortest licence duration by swapping in `min()` instead of `mean()`:

```
dog_licenses %>%
  mutate(license_duration = license_expired_date - license_issued_date) %>%
  summarise(shortest_duration = min(license_duration))
```

```
## # A tibble: 1 x 1
##   shortest_duration
##   <drtn>
## 1 1 days
```

Like `mutate()` you can also create multiple summary columns at once:

```
dog_licenses %>%
  mutate(license_duration = license_expired_date - license_issued_date) %>%
  summarise(
    avg_duration = mean(license_duration),
    shortest_duration = min(license_duration),
    longest_duration = max(license_duration)
  )
```

```
## # A tibble: 1 x 3
##   avg_duration shortest_duration longest_duration
##   <drtn>          <drtn>          <drtn>
## 1 467.321 days 1 days          2191 days
```

TODO: link to cheatsheet with list of other useful functions.

Lot's of statistical operations produce one numbers summaries and are appropriate for use with `summarise()`: `sd()`, `min()`, `max()`, `mean()`, `median()`, `quantile()` (with a single argument). Whenever you are summarizing many rows, it's a good idea to keep track of how many rows were summarized. This is so common, `dplyr` provides a special function, `n()`, that simply counts the number of rows. To add it to your summary:

```
dog_licenses %>%
  mutate(license_duration = license_expired_date - license_issued_date) %>%
  summarise(
    avg_duration = mean(license_duration),
    shortest_duration = min(license_duration),
    longest_duration = max(license_duration),
    n_licenses = n()
  )
```

```
## # A tibble: 1 x 4
##   avg_duration shortest_duration longest_duration n_licenses
##   <drtn>          <drtn>          <drtn>          <int>
## 1 467.321 days 1 days          2191 days          118542
```

Now, imagine you want this summary just for licenses issued to dogs in the Bronx. You might do something like:

Take the `dog_licenses` data, **and then**,  
 mutate to add a column called `license_duration`, **and then**  
 filter to keep rows where the borough is "Bronx", **and then**

summarise to find the mean, min and max duration along with the number of rows.

In code:

```
dog_licenses %>%
  mutate(license_duration = license_expired_date - license_issued_date) %>%
  filter(borough == "Bronx") %>%
  summarise(
    avg_duration = mean(license_duration),
    shortest_duration = min(license_duration),
    longest_duration = max(license_duration),
    n_licenses = n()
  )
```

```
## # A tibble: 1 x 4
##   avg_duration shortest_duration longest_duration n_licenses
##   <drtn>         <drtn>         <drtn>         <int>
## 1 435.9884 days 2 days          1919 days          12043
```

But how does this compare to Brooklyn? You could do the same operation again, but now for Brooklyn:

```
dog_licenses %>%
  mutate(license_duration = license_expired_date - license_issued_date) %>%
  filter(borough == "Brooklyn") %>%
  summarise(
    avg_duration = mean(license_duration),
    shortest_duration = min(license_duration),
    longest_duration = max(license_duration),
    n_licenses = n()
  )
```

```
## # A tibble: 1 x 4
##   avg_duration shortest_duration longest_duration n_licenses
##   <drtn>         <drtn>         <drtn>         <int>
## 1 465.8845 days 2 days          2191 days          29334
```

What about Queens? This kind of operation — summarising different subsets of the same data — is so common there is a much easier way to do it: combining `summarise()` with `group_by()`.

The only difference in the code, is that instead of filtering for a specific borough we'll `group_by()` the column `borough`.

```
dog_licenses %>%
  mutate(license_duration = license_expired_date - license_issued_date) %>%
  group_by(borough) %>%
  summarise(
    avg_duration = mean(license_duration),
```

```

    shortest_duration = min(license_duration),
    longest_duration = max(license_duration),
    n_licenses = n()
  )

## # A tibble: 57 x 5
##   borough      avg_duration shortest_duration longest_duration n_licenses
##   <chr>      <drtn>      <drtn>          <drtn>          <int>
## 1 ARVERNE    334.0000 days 334 days        334 days          1
## 2 Astoria    498.0000 days 239 days        757 days          2
## 3 ASTORIA    387.6667 days 366 days        405 days          3
## 4 B          347.0000 days 347 days        347 days          1
## 5 Bayside    410.0000 days 410 days        410 days          1
## 6 BELLE HARBOR 309.0000 days 309 days        309 days          1
## 7 Briarwood   540.5000 days 358 days        723 days          2
## 8 Bronx      435.9884 days  2 days        1919 days        12043
## 9 BRONX      370.3333 days 72 days         418 days          102
## 10 Brooklyn  465.8845 days  2 days        2191 days        29334
## # ... with 47 more rows

```

The `group_by()` verb doesn't perform any changes to the data except to add a signal that this data is now grouped. Subsequent operations will then happen within these groups. In the case of `summarise()` we now get one row per group, and these are all stacked together in our result.

You might have been a little surprised by the result above. I thought there were only five boroughs in New York (at least that's what the Beastie Boys told me). Notice some boroughs are represented more than once by variations in case or spelling: `Bronx`, `BRONX`. As far as `group_by()` is concerned these are distinct values of this variable. There also seem to be smaller designations than Borough in this data. You'll get a chance to try and resolve this in an exercise below.

### 6.7.1 Exercise: Dog birth months

The following code creates a new column `month_born` that holds the name of the month the licensed dog was born:

```

dog_licenses %>%
  mutate(month_born = lubridate::month(animal_birth_month, label = TRUE))

```

Use a `group_by()` step and a `summarise()` step to find the number of dogs born in each month. Which month stands out? Can you guess why?

### 6.7.2 Exercise: Order matters?

In the example above for dogs licensed in Brooklyn:

```
dog_licenses %>%
  mutate(license_duration = license_expired_date - license_issued_date) %>%
  filter(borough == "Brooklyn") %>%
  summarise(
    avg_duration = mean(license_duration),
    shortest_duration = min(license_duration),
    longest_duration = max(license_duration),
    n_licenses = n()
  )
```

the `filter()` step came after the `mutate()` step. Does this matter?

- Swap the order in the code and see if you get the same results.
- Write out how you might describe the steps. Is it obvious you can swap the filter and mutate step and get the same results?
- Which steps can't you swap the order of? Why?
- *Despite giving the same results, some orderings of the data manipulation steps will take longer to compute. Can you guess why?*

### 6.7.3 Exercise: the five boroughs

The column `neighborhood_tabulation_area` is a code for the “Neighborhood Tabulation Areas”, and has been geo-coded from the licensee’s address (as opposed to self reported). The first two characters correspond to the Borough.

This code creates a new variable called `borough_code` that contains just these two characters:

```
dog_licenses %>%
  mutate(borough_code = stringr::str_sub(neighborhood_tabulation_area, 1, 2))
```

```
## # A tibble: 118,542 x 16
##   row_number animal_name animal_gender animal_birth_month breed_name
##   <dbl> <chr> <chr> <dtm> <chr>
## 1     533 BONITA F 2013-05-01 00:00:00 Unknown
## 2     548 ROCKY M 2014-05-01 00:00:00 Labrador ~
## 3     622 BULLY M 2010-07-01 00:00:00 American ~
## 4     633 COCO M 2005-02-01 00:00:00 Labrador ~
## 5     655 SKI F 2012-09-01 00:00:00 American ~
## 6     872 CHASE M 2013-11-01 00:00:00 Shih Tzu
## 7     874 CHEWY M 2014-09-01 00:00:00 Shih Tzu
## 8     875 CHASE M 2008-08-01 00:00:00 Labrador ~
## 9     893 MILEY F 2008-07-01 00:00:00 Boxer
## 10    919 KENZI F 2010-05-01 00:00:00 Schnauzer~
## # ... with 118,532 more rows, and 11 more variables: borough <chr>,
```

```
## # zip_code <dbl>, community_district <dbl>, census_tract_2010 <dbl>,
## # neighborhood_tabulation_area <chr>, city_council_district <dbl>,
## # congressional_district <dbl>, state_senatorial_district <dbl>,
## # license_issued_date <date>, license_expired_date <date>,
## # borough_code <chr>
```

Which borough has the longest average licence duration?

## 6.8 How can I read my own tabular data into R?

### 6.8.1 What is tabular data?

Tabular data describes data that is in the form of a table: values arranged in rows each of the same length, or equivalently values arranged in columns each of the same length. Here's a small example of some tabular data:

| site_id | latitude | longitude |
|---------|----------|-----------|
| DR-1    | -49.85   | -128.57   |
| DR-3    | -47.15   | -126.72   |
| MSK-4   | -48.87   | -123.40   |

Each row records information on a site at which water measurements are taken. There are three columns: a site identification code, and the location of the site in latitude and longitude.

This is an incredibly common way of *displaying* data, but when *storing* tabular data in a file, we need a way to communicate when records and values begin and end. A very popular format for doing this is CSV. CSV, is short for comma separated values, and like the name suggests, a comma, ,, is used to separate the values for each column, while each record goes on a new line. The file is plain text, but we use the extension `.csv` to indicate that it follows the CSV format conventions.

Here's how the table above would look inside the CSV file `site.csv`:

```
site_id,latitude,longitude
DR-1,-49.85,-128.57
DR-3,-47.15,-126.72
MSK-4,-48.87,-123.4
```

In this case the first line has the column names, this is common (and recommended!), but not universal.

Why is CSV so popular?

- It's human readable. CSV isn't a special file type, it is a simple plain text file that follows some conventions. This means you don't need any special



software to look at the contents – you can open it up in anything that can examine text and take a look inside.

- It's computer readable. Because CSV files all have the same structure it's easy to write computer programs to read them. This means in almost any program designed to work with data, which is basically all the common programming languages, you'll find functions that will import CSV files. This also means it's easy to create CSV files – you can export them from Excel, write them from R, or even write one from scratch in a text editor.

If you want to look inside a CSV file in RStudio you can navigate to its location in the “Files” pane and click on its name. Selecting “View File”, will open it in the Source pane.

However, if you want to work with CSV data in R, it isn't enough to look inside the file. You need to read the contents of the file and store it in R's memory. This process is known as data import.

## 6.8.2 Importing CSV data into R

To work with data in R you need to have it in R's memory. The `read_csv()` function in the `readr` package will import a CSV file, and represent it as a tibble, if you give it the location of the CSV file. For example, to read the `site.csv` data and store it in an object called `sites`:

```
library(tidyverse)
sites <- read_csv(here("data", "site.csv"))
```

```
## Parsed with column specification:
## cols(
##   site_id = col_character(),
##   latitude = col_double(),
##   longitude = col_double()
## )
```

```
sites
```

```
## # A tibble: 3 x 3
##   site_id latitude longitude
##   <chr>      <dbl>      <dbl>
## 1 DR-1      -49.8      -129.
## 2 DR-3      -47.2      -127.
## 3 MSK-4     -48.9      -123.
```

TODO: talk about file paths, point reader to the place where file paths are talked about, or assume file is in their working directory.

Notice that `read_csv()` gave us a message about what it did: it parsed our data file and found three columns `site_id`, `latitude` and `longitude`. It also

mentions what kind of data it assumed was in each column. TODO: point to further discussion of data types.

The object `sites` is now R's representation of the data from the `site.csv` file.

### 6.8.3 Exercise: Import `visited.csv`

Use `read_csv()` to read `visited.csv` into R. How does R indicate a cell with a missing value?

```
visited <- read_csv(here("data", "visited.csv"))
```

```
## Parsed with column specification:
## cols(
##   visit_id = col_double(),
##   site_id = col_character(),
##   visit_date = col_date(format = "")
## )
```

```
visited
```

```
## # A tibble: 8 x 3
##   visit_id site_id visit_date
##     <dbl> <chr>    <date>
## 1     619 DR-1    1927-02-08
## 2     622 DR-1    1927-02-10
## 3     734 DR-3    1930-01-07
## 4     735 DR-3    1930-01-12
## 5     751 DR-3    1930-02-26
## 6     752 DR-3     NA
## 7     837 MSK-4    1932-01-14
## 8     844 DR-1    1932-03-22
```

### 6.8.4 Exercise: Import IRS tax return data for New York City

Use `read_csv()` to import the CSV file `nyc-tax-returns.csv`.

TODO: Should we introduce the “common things that go wrong” / “the most common additional arguments”, e.g. `skip`, `na`, `col_names`, `col_types`? My feeling is not now, but sometime later.

## 6.9 How can I save my results?

Say, you’ve now got a summary of the license durations by borough:

```
dog_licenses %>%
  mutate(
    license_duration = license_expired_date - license_issued_date,
    borough_code = stringr::str_sub(neighborhood_tabulation_area, 1, 2)) %>%
  group_by(borough_code) %>%
  summarise(
    avg_duration = mean(license_duration),
    shortest_duration = min(license_duration),
    longest_duration = max(license_duration),
    n_licenses = n()
  )
```

```
## # A tibble: 5 x 5
##   borough_code avg_duration shortest_duration longest_duration n_licenses
##   <chr>         <drtn>         <drtn>         <drtn>         <int>
## 1 BK           465.2768 days 2 days         2191 days         29558
## 2 BX           435.8062 days 2 days         1919 days         12050
## 3 MN           494.0276 days 1 days         2189 days         41668
## 4 QN           452.1386 days 2 days         2186 days         24420
## 5 SI           439.4875 days 4 days         2164 days         10846
```

How do you save this result for future use?

If you just need this tibble later in your code you can assign it to a variable:

```
duration_by_borough <- dog_licenses %>%
  mutate(
    license_duration = license_expired_date - license_issued_date,
    borough_code = stringr::str_sub(neighborhood_tabulation_area, 1, 2)) %>%
  group_by(borough_code) %>%
  summarise(
    avg_duration = mean(license_duration),
    shortest_duration = min(license_duration),
    longest_duration = max(license_duration),
    n_licenses = n()
  )
```

Then you can access it whenever you need it:

```
duration_by_borough

## # A tibble: 5 x 5
##   borough_code avg_duration shortest_duration longest_duration n_licenses
##   <chr>         <drtn>         <drtn>         <drtn>         <int>
## 1 BK           465.2768 days 2 days         2191 days         29558
## 2 BX           435.8062 days 2 days         1919 days         12050
## 3 MN           494.0276 days 1 days         2189 days         41668
## 4 QN           452.1386 days 2 days         2186 days         24420
```

```
## 5 SI          439.4875 days 4 days          2164 days          10846
```

This keeps our result around in memory, but often you also want to preserve the data in a file on disk. There are two common choices for format: CSV and RDS.

You’ve already seen CSV files. Saving your results in this format gives you all the benefits of that format: plain text files easily shared and opened. You can save a tibble as a CSV file with the `readr` function `write_csv()`, where all you need to specify is the path:

```
duration_by_borough %>% write_csv("duration-by-borough.csv")
```

You can then read this file in any project or R session with `read_csv()`:

```
duration_by_borough <- read_csv("duration-by-borough.csv")
```

```
## Parsed with column specification:
## cols(
##   borough_code = col_character(),
##   avg_duration = col_double(),
##   shortest_duration = col_double(),
##   longest_duration = col_double(),
##   n_licenses = col_double()
## )
```

RDS files are a special R format. They are binary files as opposed to plain text files, which means you can’t just open them up and look inside. If you are sharing them you’d also need your collaborators to have R. These are downsides, but the advantage of this format is that it can be much quicker to load, it will preserve special R data types (for example factors, or nested structures), and can save any R object not just tabular data structures.

To save the tibble as an RDS, use `write_rds()`:

```
duration_by_borough %>% write_rds("duration-by-borough.rds")
```

To read it back in, use `read_rds()`:

```
duration_by_borough <- read_rds("duration-by-borough.rds")
```

Often you’ll save your data in both formats to make sure you get the best of both worlds. You’ll talk more about where to save your data in [TODO: add link](#).

### 6.9.1 Exercise

Save the tibble with our extra columns:

```
dog_licenses %>%  
  mutate(  
    license_duration = license_expired_date - license_issued_date,  
    borough_code = stringr::str_sub(neighborhood_tabulation_area, 1, 2))
```

into a CSV file called `dog-licenses-extra.csv`.



# Chapter 7

## Publishing

### 7.1 Questions

- How can I share my work on the web?

### 7.2 Why should I share my work on the internet?

A key part of any project is communicating what you have learned. You've already seen how to create documents to communicate your work with R Markdown, and how to host whole projects on GitHub. This chapter is about sharing your work through a webpage on the internet. Some advantages of sharing your work on a webpage include:

- It's easy for your visitors — they just need to click on a link and see your work. They don't need to know anything about R Markdown, HTML, or GitHub.
- It provides a visually friendly and customizable landing point for people interested in your project. You can easily point them to the GitHub repo if they want more details.
- It's easy for you to make updates and those updates are immediately available to any visitors — you don't have to re-send anyone any files.

By the end of this chapter, you'll be able to take an R Markdown document that lives on your computer and share it with the world through a link to your own webpage.

## 7.3 What does it take to get a webpage online?

To understand what it takes to get a webpage online, it helps to understand roughly what happens when you point your browser at a web address. When you point your browser at an address, for example `https://merely-useful.github.io/r-publishing.html`, the following things happen:

1. The name `https://merely-useful.github.io` is converted to an address of a computer that has the relevant files for this website. This computer is also known as the website host or server.
2. The browser requests the file `r-publishing.html` from the host.
3. The browser reads the contents of `r-publishing.html` and displays it for you in the browser window.

In reverse order this process also describes what you need to do to get your own website online:

1. You need an HTML file that describes what people should see on your page.
2. You need to host the HTML file on a computer on the internet.
3. You need a way to associate a URL with the address of your host.

In this chapter, you'll learn a process for getting your work online that leverages what you already know — creating HTML files using R Markdown (#1 above) and how to host your work in a GitHub repository (#2). The third step will be handled by a GitHub service called GitHub Pages. By following the conventions that GitHub Pages expects, you'll be able to make a webpage for any of your repositories available at: `http://{{your_username}}.github.io/{{repo_name}}`.

## 7.4 How do I get my work on the web?

### 7.4.1 A starting point

TODO: revisit this later when more content is fleshed out. Maybe there will be a repo we can rely on all learners having that we can start from.

In practice you'll probably start thinking about a website once you've already done a lot of work on a project — your project will already have some analysis documented in R Markdown, be in version control, and hosted on GitHub. However, so that we can work with a specific example, you'll set up a project in this section that is less developed than where your project might be when you start thinking about making a website.

At the end of this section, you should have an RStudio project with an example report in `report.Rmd`. This project should also be on



GitHub at `https://github.com/{{your_username}}/sharing-work`, where `{{your_username}}` should be substituted with your GitHub username, e.g. mine is at `https://github.com/cwickham/sharing-work`.

Let's start by creating a new project called, `sharing-work` by running the following in the RStudio console:

```
usethis::create_project("sharing-work")
```

Once the project opens, set it up to use version control by running the next code in the Console:

```
usethis::use_git()
```

Then add it as a repository on GitHub:

```
usethis::use_github()
```

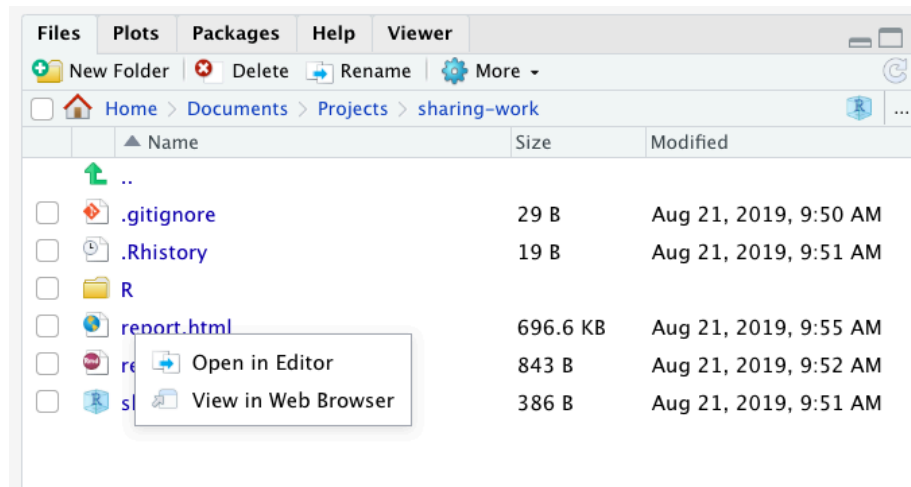
Then so we have a report to work with, create a new R Markdown file (“File -> New File -> R Markdown”), making sure to leave the “Default Output Format” as HTML. Save this new file as `report.Rmd`.

Commit these changes and push your repository to GitHub. This is our starting point.

TODO: Would it be better to get to this point by getting learners to fork a repo, then “New project -> From version control” in RStudio? Except forking isn't in the plan for the Version Control section.

### 7.4.2 HTML files

If you Knit `report.Rmd` you'll get `report.html`, an HTML document, because in the header of `report.Rmd` output is set to `html_document`. In the Files pane in RStudio if you click on `report.html`, you'll get two options: Open in Editor, or Display in Web Browser.



HTML is the language of webpages. If you “Open in Editor” you will see the contents of the file — it’s plain text, but follows the conventions of the HTML language. If you “View in Web Browser” your browser will read, interpret and display the HTML for you. When you “View in Web Browser” you may notice the address bar in your browser looks something like:

```
file:///Users/wickhamc/Documents/Projects/sharing-work/docs/report.html
```

The `file://` is a signal that this particular address is to a file living on your computer locally. You couldn’t give this address to someone else and expect it to work, because they don’t have this file, let alone the same directory structure.

The HTML produced by R Markdown is completely self-contained, the browser needs no additional files to display the page as you see it now. So, you could email the file `report.html`, and your recipients could open it their browser and see the same result. However, rather than sending the file, our goal will be to put this HTML file on the web so you can share a link. You’ll start by having this file accessible at the link `https://{{your_username}}.github.io/sharing-work/report.html`, then learn how to have it displayed with the shorter link `https://{{your_username}}.github.io/sharing-w`

### 7.4.3 Setting up your repo to have a web page

Our goal is to get the report that is currently living in `report.html` displayed when a visitor heads to `https://{{your_username}}.github.io/sharing-work`

You’ll set up GitHub Pages to look for your HTML files in the `docs` directory. Your first step will be to make this directory and put our report file inside. Create `docs/` with:

```
usethis::use_directory("docs")
```

Move `report.Rmd` and `report.html` into this directory. You can do this within the Files pane by checking the files and “More -> Move...”. Or you can do it with code on the Console:

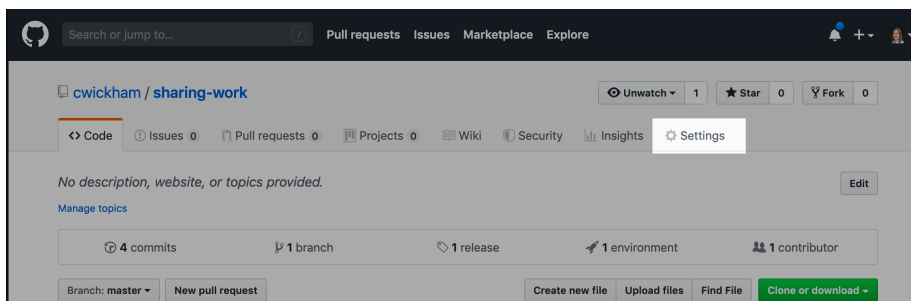
```
file.rename("report.Rmd", "docs/report.Rmd")
file.rename("report.html", "docs/report.html")
```

Commit and push your changes. Your repo should now have a structure like:

```
docs
  report.Rmd
  report.html
sharing-work.Rproj
```

Now you will activate GitHub Pages, so that GitHub will know to deliver your files when visitors head to: `https://{{your_username}}.github.io/sharing-work`

Visit your GitHub Repository and head to the “Settings” tab.



Scroll to the “GitHub Pages” section. Activate GitHub pages, with source set to “master branch /docs folder”. You should see a message that your site is now live at: `https://{{your_username}}.github.io/sharing-work`

## GitHub Pages

GitHub Pages is designed to host your personal, organization, or project pages from a GitHub repository.

Your site is ready to be published at <https://cwickham.github.io/sharing-work/>.

### Source

Your GitHub Pages site is currently being built from the `/docs` folder in the `master` branch. [Learn more.](#)

master branch /docs folder ▾

Try visiting: `https://{{your_username}}.github.io/sharing-work/report.html`. You should see your report. Congratulations you have a webpage!

#### 7.4.4 Getting a default page to display when people visit the project site

You could send people the link, `https://{{your_username}}.github.io/sharing-work/report.html`, but it is often nicer to send them the shorter version without the file name: `https://{{your_username}}.github.io/sharing-work`. You can try this now, but it won't work — you'll see a message in your browser like: "404- File not found". This shorter URL points to a directory as opposed to a file. By default, when a server receives a request for a directory, it looks for a file to display with a default name — usually `index.html`. In your case there is no file called `index.html` so there is nothing to display.

If you would like the contents of `report.html` to be displayed as the homepage of your project, then rename `report.Rmd` to `index.Rmd`. You'll then need to regenerate the HTML file, remove `index.html`, commit, and push your changes. For work that is communicated easily in one page, this would be a good option.

Alternatively, you might have a different page as the default page — one that summarizes the project and then links to other more detailed pages. You'll see how to do this over the next few sections. To get started create a new R Markdown document. Delete **all** the contents of the file, and then copy and paste in the following:

```
---
title: "Sharing work on a webpage"
author: "Me"
output: html_document
---
```

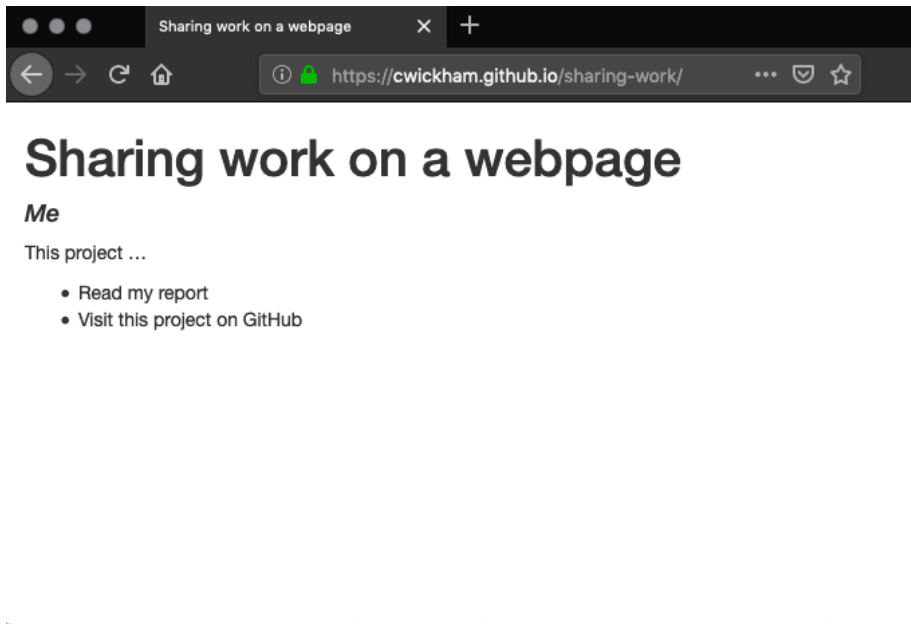
This project ...

- \* Read my report
- \* Visit this project on GitHub

Save the file as `index.Rmd`, knit it, commit, and push your changes. Your repo should now look like:

```
docs
  index.Rmd
  index.html
  report.Rmd
  report.html
sharing-work.Rproj
```

Now when you visit `https://{{your_username}}.github.io/sharing-work` you should see:



#### 7.4.5 What does it take to get your work on a webpage?

To sum up the process above, in its most minimal form, to have a webpage at `https://{your_username}.github.io/{repo_name}`, your repo at `https://github.com/{your_username}/{repo_name}` needs to:

1. have an `index.html` file in the `docs` directory (probably generated from `index.Rmd` in the same location), and
2. have GitHub Pages activated in repository settings with source set to “master branch /docs folder”.

Be aware that everything inside the `docs` folder is now public, even if your repository is private.

You might have noticed this book lives at a GitHub Pages URL without a repo name — there is nothing after the `.io` in `https://merely-useful.github.io`. This is known as a user or organization site (`merely-useful` is an organization rather than a user, so this is an organization site). There are a few differences between what you’ve learnt so far and the process of setting up a user site at `https://{your_username}.github.io`. First, you need to name your repository in a specific way — it must be called `{your_username}.github.io`. Second, user sites don’t use the `docs/` folder — you put your HTML files at the top level in the repo. And third, you don’t have to change any settings with user sites — GitHub will recognize the repo name and automatically serve it at `https://{your_username}.github.io`.

### 7.4.6 Exercise: Customize `index.Rmd`

- Edit `index.Rmd` to have your name as the author.
- Knit `index.Rmd` to verify your changes, then commit and push them.
- Visit `https://{{your_username}}.github.io/sharing-work` to check the updated site.

*This is the workflow for making changes to your webpage. Make edits locally, and Knit to check them. Then commit and push to make those changes visible on the web.*

## 7.5 How do I link to other pages, files or images?

### 7.5.1 Linking to other pages

To create a link to another page in markdown file you use the syntax:

```
[text to display](url)
```

Once Knit to HTML, only `text to display` will be visible, and clicking on the text will take a viewer to `url`. For example, to add a link to my GitHub repo I might add the following line to `index.Rmd`:

```
Visit [my github repo](https://www.github.com/cwickham/sharing-work)
```

Which when Knit to HTML renders like:

Visit my github repo

This is an example of an **absolute** URL. Just like when you specify file paths on your own computer, URLs can be both absolute and relative. An absolute URL describes a file location starting from and including the domain name. For instance, the absolute URL that points to `report.html` in my repo is `https://cwickham.github.io/sharing-work/report.html`.

Relative URLs are *relative* to the current HTML file. So, for instance if you are viewing `https://cwickham.github.io/sharing-work/index.html`, a relative link to my `report.html` would be `report.html` since this file is at the same level as `index.html` in my website structure. For pages created using this GitHub Pages workflow, your website structure is the same as the file structure in your `docs/` folder.

To add a link to `report.html` in `index.Rmd` I would add a line like:

```
See the [full report](report.html)
```

You should use relative URLs to reference any of **your** files (i.e. those in `docs/`). That way if you ever rename your repository, move it, or use a different hosting platform, your links will all work without changes. You must use absolute links for files that reside elsewhere on the internet.

### 7.5.2 Exercise: Relative links

Imagine your `docs/` folder had the following structure:

```
docs
  index.Rmd
  index.html
  diagrams
    workflow.png
  |   reports
  |     jan.Rmd
  |     jan.html
  |     feb.Rmd
  |     feb.html
  |   sharing-work.Rproj
```

Using a relative URL, how would you refer to:

- `jan.html` from `index.html`?
- `feb.html` from `jan.html`?
- `workflow.png` from `index.html`?

### 7.5.3 Exercise: Add links to `index.html`

Add to `index.Rmd`:

- a link to `report.html` using a relative link, and
- a link your GitHub repository using an absolute link.

### 7.5.4 Linking to sections within a page

URLs can also refer to places inside the current page, most usually to another section. In R Markdown you've seen you create headings using `#`. For example, an Appendix subsection might be:

```
## Appendix
```

If we want to link to this section from elsewhere, you prefix the section name with `#` in the URL. For example, if this section is in `report.Rmd` and you want to link to it elsewhere in `report.Rmd`, you could use:

See more details in the [\[appendix\]\(#appendix\)](#)

The URL `#appendix` is interpreted as the heading with ID `appendix` in the current page. R Markdown creates IDs for all sections (and subsections) automatically, by converting to lower case and replacing spaces with dashes (-). But, you can explicitly set IDs too, by adding the ID with the `#` inside curly braces after the section heading. For instance you might prefer the shorter `appen` ID. You need to set it where the heading occurs:

```
## Appendix {#appen}
```

Then you can link to it using this shorter ID elsewhere:

See more details in the [\[appendix\]\(#appen\)](#)

You can also use this strategy to link to sections in other pages by including the relative URL first. For instance, to refer to this “Appendix” section from `index.html` you could include in `index.Rmd`:

Some gory details of the analysis can also be found in the [\[Appendix of the report\]\(report.html#appen\)](#)

### 7.5.5 Exercise: Add and link to a section in `report.Rmd`

Add a new section to `report.Rmd` and include a link to it in `index.Rmd`. You’ll need to Knit both `report.Rmd` and `index.Rmd`, and commit and push the HTML files to check your work.

### 7.5.6 Including images

Most of your images will likely be generated by R chunks in your R Markdown files and be included automatically. If you want to include display other images, you use the same syntax you saw in the Markdown section of the Reproducibility chapter. That is, in your R Markdown file you’ll include the image with something like:

```
![Image caption here.](path/to/image/file.png)
```

However, the `path/to/image/file.png` should be a relative URL pointing at an image in your `docs/` directory. For example, if you had an image, `me.png`, inside an `images` directory inside your `docs` folder:

```
docs
  index.Rmd
  index.html
  images
    me.png
```



```
report.Rmd
report.html
sharing-work.Rproj
```

You could include it in `index.html` by adding to `index.Rmd` the line:

```
![A picture of me](images/me.png)
```

Notice the syntax is very similar to adding a link to `me.png`:

```
[A picture of me](images/me.png)
```

*Including* the image displays the image inside at the appropriate place in the current page, *linking* to the image requires a viewer to click the link to see the image.

Your image will be included at full size, but you might find it too large. You can additionally specify some attributes for the image in curly braces immediately following the link. For instance, use the `width` attribute to set the image width, either in pixels:

```
[A picture of me](images/me.png){width=50} # default unit is px
```

Or as a percentage:

```
[A picture of me](images/me.png){width=50%}
```

### 7.5.7 Exercise: Add an image to `index.html`

Include an image in your `index.Rmd`. (If you need an image to include you could always build your own version of an Octocat, GitHub’s mascot).

Don’t forget, you’ll need to commit your re-Knit `index.html` and your image.

## 7.6 Exercise: Add a website to an existing project

Add a website to one of your existing project repositories. You’ll need to complete the following steps:

- Create a `docs/` directory in your project.
- Add an `index.Rmd` R Markdown document to the `docs/` folder, knit it to HTML to produce `index.html`.
- Commit these changes to git, and push to GitHub.
- Activate GitHub Pages in the repository settings with source set to “master branch /docs folder”.
- Visit the site to check it is working.



## Part II

# Novice Python Material



## Chapter 8

# Introduction

FIXME: general introduction.

### 8.1 Who are these lessons for?

FIXME: personas

#### 8.1.1 Summary

FIXME: elevator pitch

#### 8.1.2 Prerequisites

FIXME: prerequisites

### 8.2 What does “done” look like?

FIXME: end goal

### 8.3 What will we use as running examples?

FIXME: introduce running examples



## Chapter 9

# Development

### 9.1 Questions

- What are functions and how can I make my own?
- How can I make my programs tell me that something has gone wrong?
- How can I ask for help online?
- What are packages and how do I install them?

### 9.2 Objectives

- Describe why functions are used and are major components of programming.
- Create a function that takes both positional arguments and keyword arguments.
- Use `assert` to check that a function input parameter is within a certain range of values.
- Formulate a question with a minimal reproducible example.
- Describe why there are packages and why not all code is included with the default Python installation.
- Install and use a package.

### 9.3 Functions

Functions are like recipes. You give a few ingredients as input to a function, and it will generate an output based on these ingredients. Just as when following a recipe, both the ingredients and the instructions will influence the final result.

In Python, the inputs to a function are not called ingredients, but rather arguments, and the output is referred to as the return value of the function. A function does not technically need to return a value, but often does so. Functions facilitate reusing chunks of code in a way that is more readable and reproducible than cutting and pasting several lines of code. E.g. if our data analysis code is broken down into functions, we could readily use it with many different data sets by changing the input data path, but leaving the rest of the code the same.

Well chosen function names also clarifies the flow of analysis. For example, imagine that you open a file with the following lines of code within it.

```
images = read_in_images(file_paths)
gray_images = convert_to_grayscale(images)
brightest_image = find_brightest_image(gray_images)
```

Just by looking at the function names, it is clear what this code is intended for and its main flow of operations is immediately visible. Inside each of these functions there might be 10-20 lines of code, so if we would not have modularized the code into separate functions with well chosen names, it would take longer to understand its overall purpose since there would be 30-60 lines of code to read instead of just three.

### 9.3.1 Creating functions

There are many useful functions already built into Python, and the ability to create your own allows you to string together any sequence of operations in ways that are tailored to your workflow. The `def` keyword lets us define a function with a name of our choice and an arbitrary number of input parameters.

```
def sum_two_numbers(num1, num2):
    return num1 + num2
```

This function accepts two input parameters, `num1` and `num2`, and returns their sum. Just as with variable names, function names are preferably written in **snake\_case** (see the style guide for details), and avoid existing Python keywords and built-in names (a list of these is available [\[here\]](#)<sup>[so-keywords-builtins]</sup>, but instead of memorizing that list you can type the desired name into the Python interpreter to find out if it already exists).

To execute the operations listed in the function, we can call the function and pass the two numbers we want to add as the arguments to the function.

```
sum_two_numbers(2, 5)
```

```
## 7
```

The returned value can be assigned to a variable:



```
number_sum = sum_two_numbers(2, 5)
number_sum
```

```
## 7
```

A more versatile function could add any amount of numbers together and return their sum:

```
def sum_all_numbers(list_of_numbers):
    number_sum = 0
    for number in list_of_numbers:
        number_sum += number
    return number_sum
```

```
sum_all_numbers([1, 2, 3,])
```

```
## 6
```

A function can also return multiple outputs, e.g. we can return the number of elements in addition to their sum:

```
def sum_and_len_all_numbers(list_of_numbers):
    number_sum = 0
    number_len = 0
    for number in list_of_numbers:
        number_sum += number
        number_len += 1
    return number_sum, number_len
```

To capture the output, we can either assign to a single name, a tuple, or assign both at the same time to different variables.

```
sum_and_len_of_numbers = sum_and_len_all_numbers([1, 2, 3])
sum_and_len_of_numbers
```

```
## (6, 3)
```

```
sum_of_numbers, len_of_numbers = sum_and_len_all_numbers([1, 2, 3])
sum_of_numbers
```

```
## 6
```

```
len_of_numbers
```

```
## 3
```

Note that when we defined the function `sum_two_numbers()`, we referred to `num1` and `num2` as *parameters*, while we refer to the numbers we pass to the function call (2 and 5 above) as *arguments*. Although this might sound confusing at first, it is a standard followed by many programming languages so it is useful to get accustomed to this terminology.

### 9.3.1.1 Function composition

Combining functions is referred to as function composition. This practice allows us to write functions that perform one specific task and then combine them for more complicated tasks, which makes code more readable and easier to debug. By composing a function from the built-in `len()` and `sum()` functions, we can create a more succinct and easier to read version of our previous function `sum_and_len_all_numbers()`.

```
def sum_and_len(list_of_numbers):  
    return sum(list_of_numbers), len(list_of_numbers)
```

### 9.3.1.2 Positional and keyword arguments

Up until now, our function calls have included just enough arguments to assign one to each of the function parameters. The assignment has been based on the position of the arguments in the function call and therefore these are called positional arguments. We could be more explicit and include the parameter name in the assignment.

```
sum_two_numbers(num1=2, num2=5)
```

```
## 7
```

The arguments are now referred to as keyword arguments and has the advantage that they can be specified in any order.

```
sum_two_numbers(num2=5, num1=2)
```

```
## 7
```

### 9.3.1.3 Defining default values

Above, we have specified a value for every argument in each function call. This is manageable when functions have few parameters, but it can get tedious for functions with many parameters. Defining default values for select parameters can facilitate working with complex functions by reducing the number of arguments that need to be defined when calling the function and guide users to good default choices for the parameters without requiring in-depth knowledge of each parameter. As an example, we can modify `sum_two_numbers()` to optionally return the two input arguments.

```
def sum_two_numbers(num1, num2, return_input=False):  
    if return_input:  
        return num1, num2, num1 + num2  
    else:  
        return num1 + num2
```

By default the function will work just as previously.

```
sum_two_numbers(2, 5)
```

```
## 7
```

But we now also have the option to return the input numbers.

```
sum_two_numbers(2, 5, return_input=True)
```

```
## (2, 5, 7)
```

Since the arguments are given in the same order as in the function definition, we could have left out the `return_input=` part and just written `True` in the third position.

#### 9.3.1.4 Function documentation

Functions might appear self-explanatory when they are being written, but it is essential that there is proper documentation describing what the function does and what types of arguments should be in the input. This helps other people who are reading your code and also your future self that will be reusing these functions.

In Python, a function is documented in its docstring, which is a multiline Python string immediately following the function definition. It is surrounded by triple quotes (either single or double) and can look like this.

```
def sum_two_numbers(num1, num2):  
    '''  
        Add two numbers  
  
        Parameters  
        -----  
        num1: int, float  
            The first number to be added.  
        num2: int, float  
            The second number to be added.  
  
        Returns  
        -----  
        int, float  
            Sum of the two numbers.  
    '''  
  
    return num1 + num2
```

The above docstring convention is referred to as the `numpy` docstring format. There are other conventions, but here we recommended using the `numpy` doc-

string format since it is easy to read for complex functions with many arguments and used by many integral data science Python packages. It is described in great detail in the `numpy` and `pandas` documentation.

Docstrings constitute the text that show up in the function help message, so it is important that these are well-written and helpful for the reader.

```
help(sum_two_numbers)
```

```
## Help on function sum_two_numbers in module __main__:
##
## sum_two_numbers(num1, num2)
##     Add two numbers
##
##     Parameters
##     -----
##     num1: int, float
##         The first number to be added.
##     num2: int, float
##         The second number to be added.
##
##     Returns
##     -----
##     int, float
##         Sum of the two numbers.
```

This is the same text that is displayed in the Spyder help pane, where it is rendered as markdown for rich display with typefaces, headings, etc.

Spyder includes a convenient function to automatically generate docstring templates. Once you have written the function signature and typed out the first triple quote for the docstring, there will be a small pop-up window that reads “Generate docstring”. Click it or press enter and a `numpy` docstring template will be created based on the parameters in the function signature.

## 9.4 How to make programs indicate that something has gone wrong?

When performing programmatic data analysis, we need to both ensure that our code runs correctly and that it carries out the intended tasks. The Python interpreter will help us with the first part; if a part of the code is not valid, an error will be raised with a message that helps us trace back what part of the code is not correct.

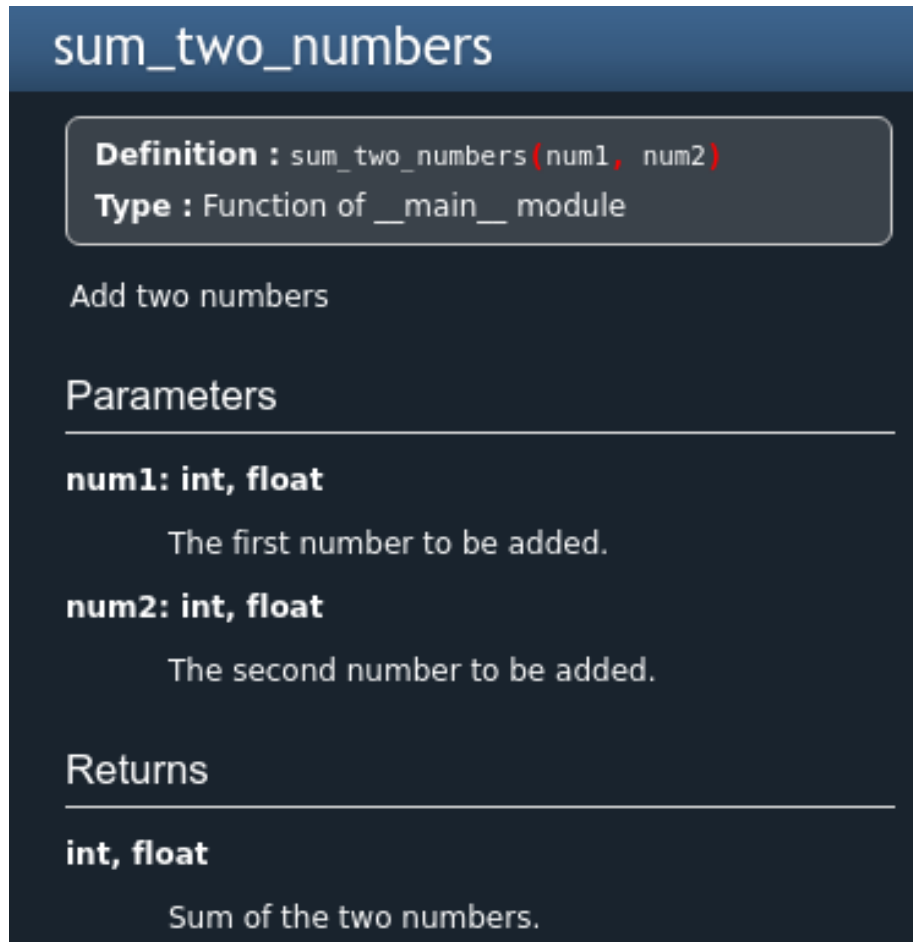


Figure 9.1: Rich rendering of the docstring in Spyder

### 9.4.1 Ducktyping - relying on Python to detect unexpected behavior

In our function `sum_two_numbers()`, we could explicitly check that the input arguments are numerical. However, explicitly checking the type of each input parameter quickly becomes tedious and can make functions less readable, especially as they grow more complex. An alternative approach is to try to perform the intended operations on the input parameters and rely on that Python will raise an error if they are of the wrong type.

```
sum_two_numbers(5, 'six')
```

```
## Error in py_call_impl(callable, dots$args, dots$keywords): TypeError: unsupported op
##
## Detailed traceback:
##   File "<string>", line 1, in <module>
##   File "<string>", line 18, in sum_two_numbers
```

The raised error contains a helpful message that alerts us to what went wrong; `string` objects cannot be added to `int` objects in Python. If the input variables behave correctly then they probably are of the correct types and no explicit checking is needed. This approach is often referred to as “ducktyping” because it makes assumptions on the type of variable based on its behavior, just like the saying

If it looks like a duck, swims like a duck, and quacks like a duck,  
then it probably is a duck.

### 9.4.2 Assertions - explicitly checking for unexpected behavior

Ducktyping is useful to catch anything that raises an error in Python, but sometimes we might want to stop the code execution for reasons other than a technical Python error. Examples of this include when you have performed a specific operation that you know should give output of a certain shape and to check if a variable is within an allowed range. The `assert` statement allows us to check if a condition is `True` and stop code execution if it is not.

```
assert 1 == 0
```

```
## Error in py_call_impl(callable, dots$args, dots$keywords): AssertionError:
##
## Detailed traceback:
##   File "<string>", line 1, in <module>
print('This is only printed if the assertion above is "True"')

## This is only printed if the assertion above is "True"
```

Since the condition above is `False`, an `AssertionError` is raised and code execution halts. If the assertion is `True`, there is no output and the next line of code is executed (if there is one).

```
assert 0 == 0
print('This is only printed if the assertion above is "True"')
```

```
## This is only printed if the assertion above is "True"
```

As mentioned above, this is useful if we know that a variable should look a certain way, since we can assert if this is the case and guard ourselves from errors that originate early in the pipeline but could give rise to more cryptic errors that are difficult to troubleshoot later in the pipeline.

It is often helpful to add a clarifying message to the assertion statement, especially as assertions become more complex.

```
x = 1
assert x == 0, 'x is not 0'
```

```
## Error in py_call_impl(callable, dots$args, dots$keywords): AssertionError: x is not 0
##
## Detailed traceback:
##   File "<string>", line 1, in <module>
```

Note that similar functionality can be achieved by explicitly raising an error within a `try` and `except` block, or within a conditional statement using `if`, `elif`, and `else`, but `assert` is a simple and readable way of allowing for manual error checking.

## 9.5 Packages

### 9.5.1 Installing Python

There are several ways of installing Python on your system. One of the most robust and cross-platform compatible is to install the Anaconda Python distribution, which is available for Linux, macOS and Windows. Choose to download the Python 3 installer unless you need to work with Python 2 for a specific reason.

### 9.5.2 Using Python

When the Anaconda installation has finished, Python is accessible by running `python` in a terminal (on Windows, use the `Anaconda Prompt`). Anaconda also includes graphical interfaces for interacting with Python, which can be invoked

by running `spyder` or `jupyter-lab` from the command line, these are covered more in detail elsewhere in this book.

### 9.5.3 What is a package?

Certain functionality that is considered essential for the Python programming language is available wherever Python is installed. Other highly useful, but often more domain specific functionality can be accessed separately in the form of Python packages. A package is essentially a few Python scripts in a specific directory structure coupled with installation instructions for the computer.

Note that you will sometimes hear packages refers to as “modules”. The two words are often used interchangeably. Technically, a package is a folder that contains modules (scripts), which in turn contains functions (code).

People around the world have created packages for Python and made them freely available for others to use resulting in one of the richest package ecosystems for any programming language, with packages for web design, prose writing, game development, and data science (to name just a few). Since there are so many packages available, it is not feasible to include all of them with the default Python installation (it would be as if your new phone came with every single app from the app/playstore preinstalled). Instead, Python packages can be downloaded from central repositories online and installed as needed. The two main repositories are the Python Package Index (PyPi) and the Anaconda package repository. Instead of navigating to these sites with a web browser, downloading the desired packages and installing them manually, there are so called package managers that automate these processes. The package manager for PyPi is called `pip` and the package manager for Anaconda is called `conda`. To install a package either of these can be used, below is an example of what to type on the command line to install the **numerical python** package `numpy`.

```
conda install numpy
```

```
pip install numpy
```

Uninstalling a package is equally simple

```
conda remove numpy
```

```
pip uninstall numpy
```

Since we have downloaded the Anaconda Python distribution, we will predominantly be using the `conda` package manager. Installing packages with both interchangeably works, but it is recommended to stick to one as much as possible. The Anaconda team has already bundled many commonly used packages together with their Python installer, including most of the common data science packages, such as `numpy` and `pandas`.



Now that you know how to access many of the world's best data science packages right from your terminal, let's see how we can use them!

Pro tip: Some packages are not available in the default Anaconda repositories. User contributed packaged are available in Anaconda "channels", use `anaconda search -t conda <package name>`, to find a channel with the desired package. To install this package, use `conda install -c <channel name> <package name>`. The conda forge channel channel has many of the packages not in the default repositories.

### 9.5.4 Importing packages

Any installed package can be accessed by typing `import <package_name>` in Python, e.g.

```
import numpy
```

After importing a Python package, we can access any of its functions, by first writing the followed by a period and then the function name, e.g.

```
numpy.mean([1, 2, 3])
```

```
## 2.0
```

You can think of this as navigating to the numpy menu in a GUI software, and then clicking the function you want. You don't need to recall every function name by hard, pressing the TAB key after the period, will bring up all available function and intelligently filter them as you type out more letters, try it! When you start getting familiar with typing function names, you will notice that this is often faster than looking for functions in menus.

It is common to give packages shorter nicknames, which are faster to type. This is not necessary, but can save work in long files and make code less verbose so that it is easier to read:

```
import numpy as np
```

```
np.mean([1, 2, 3])
```

```
## 2.0
```

We could also import the mean function directly.

```
from numpy import mean
```

```
mean([1, 2, 3])
```

```
## 2.0
```

And even give it a nickname.

```
from numpy import mean as mn  
  
mn([1, 2, 3])
```

## 2.0

Which of these you use is up to you, but it is common to follow the conventions established by the library’s authors, which for numpy is `import numpy as np` and use functions via `np.<function_name>`. One thing to avoid is to import everything from a package, e.g. `from numpy import *`. If this is done with every package, it is almost guaranteed that the same function name will be available from more than one package and it will be difficult to keep track of it if you are using the `mean` function from `numpy` or from another package that you also imported everything from.

Subpackages and their functions can be imported via the dot syntax.

```
from numpy.fft import fftfreq
```

Packages and subpackages might sound complicated to keep apart. It can be helpful to understand that they are folders and files in specific directory structure. Considering only the `numpy` packages we have mentioned so far, the directory structure would look like this.

```
numpy-folder/  
  script-with-the-mean-definition.py  
  fft-folder/  
    script-with-the-fftfreq-definition.py
```

### 9.5.5 Installing packages

After Anaconda has been installed on your system, you can use the command line `conda` package manager or the GUI-driven `anaconda-navigator` to install Python packages. For comprehensive instructions on both of these, refer to the official documentation. Brief step-by-step instructions to get up and running with `conda` follow.

1. To install a new Python package from the Anaconda repositories, simply run `conda install <package name>` in a terminal. You can also use the `pip` package manager, but it will be easier to keep track of packages by sticking to one installation method.
2. Some packages are not available in the default Anaconda repositories. User contributed packaged are available in Anaconda “channels”, use `anaconda search -t conda <package name>`, to find a channel with the desired package. To install this package, use `conda install -c`

`<channel name> <package name>`. The conda forge channel has many of the packages not in the default repositories.

## 9.6 How to get help online

When reading the built-in Python help is not sufficient, there are several online resources that can be helpful. One of the most commonly used resources for data science is the Stack Exchange network which offers Q&A sites both for programming related topics via Stack Overflow, as well as statistics and machine learning via Cross Validated.

If nothing relevant can be found after searching the many existing questions and answers on these sites, it is appropriate to ask a new question! Stack Overflow has detailed instructions on how to create a minimal reproducible example when asking a question to increase the chances that the question receives a specific and helpful answer. The key principles listed on the website recommend that an answer follows these guidelines:

- Minimal – Use as little code as possible that still produces the same problem
- Complete – Provide all parts someone else needs to reproduce your problem in the question itself
- Reproducible – Test the code you’re about to provide to make sure it reproduces the problem

An additional benefit of reducing the problem into this format is that you might discover the error in the process! This process helps narrowing down exactly which region of the code is failing and the act of explaining the problem often reveals a solution before any replies have come in (commonly referred to as rubber duck debugging).

## 9.7 Key Points

- Functions allow partitioning of code into segments that perform specific tasks. They make code easier to write, read, troubleshoot, and reuse.
- Functions can have default values for certain input parameters, which allows these parameters not to be explicitly set every time a function is called.
- Ducktyping relies on the Python interpreter to raise meaningful errors when an illegal operation tries to execute.
- Assert statements can be used to explicitly check the truth value of a condition and halt code execution if needed.
- When asking for help online, it is critical to include a minimal reproducible example so that the question can receive suitable replies.

- Packages makes it easy to reuse and share code in an organized and automated manner. These can be installed with package managers.

## 9.8 Exercises

1. Write a function that returns the mean of exactly two numbers. Start with the code we used to create the function `sum_two_numbers()` above.
2. Write a function that returns the mean of any amount of numbers. Start with the code we used to create the function `sum_and_len_all_numbers()` above.
3. Write a function that takes a string as its input and returns a tuple with the first and last character of the string.
4.
  - a. Write a function that takes a string as its input and can return any single character. Which character is returned should be determined by an integer passed to a parameter called `character_index` which has the default value of 0.
  - b. Expand on the previous function by including an assert statement that checks if the integer given to the `character_index` parameter is within the length of the string.
  - c. Describe the pros and cons of the ducktyping approach in 4a and the explicit assert approach in 4b.

**Part III**

**RSE Material**



## Chapter 10

# RSE Introduction

As research becomes more computing intensive, researchers need more computing skills so that:

- other people (including your future self) can re-do your analyses;
- you and the people using your results can be confident that they’re correct; and
- re-using your software is easier than rewriting it.

Most books and courses about software engineering are aimed at product development, but research has different aims and needs. A research programmer’s goal is to answer a question; she might build software in order to do that, but the software is only a means to an end.

But just as some astronomers spend their entire careers designing better telescopes, some researchers choose to spend their time building software that will primarily be used by their colleagues. People who do this may be called research software engineers (RSEs) or data engineers, and the aim of these lessons is to help you get ready for these roles—to go from writing code on your own, for your own use, to working in a small team creating tools to help your entire field advance.

One of the many challenges you will face is to find the appropriate mix of tools and methods for each problem you have to solve. If you want to reformat a handful of text files so that your program can read them in, you shouldn’t bother writing a comprehensive test suite or setting up automated builds. On the other hand, if you *don’t* do this, and that “handful of text files” turns into a pile, and then a mountain, you will quickly reach a point where you wish you had. We hope this training will help you understand what challenges have already been solved and where to find those solutions so that when you need them, you’ll be able to find them.

## 10.1 Who are these lessons for?

**Amira** completed a Master’s in library science five years ago, and has worked since then for a small NGO. She did some statistics during her degree, and has learned some R and Python by doing data science courses online, but has no formal training in programming. Amira would like to tidy up the scripts, data sets, and reports she has created in order to share them with her colleagues. These lessons will show her how to do this and what “done” looks like.

**Jun** completed an Insight Data Science fellowship last year after doing a PhD in Geology, and now works for a company that does forensic audits. He has used a variety of machine learning and visualization software, and has made a few small contributions to a couple of open source R packages. He would now like to make his own code available to others; this guide will show him how such projects should be organized.

**Sami** learned a fair bit of numerical programming while doing a BSc in applied math, then started working for the university’s supercomputing center. Over the past few years, the kinds of applications they are being asked to support have shifted from fluid dynamics to data analysis. This guide will teach them how to build and run data pipelines so that they can teach those skills to their users.

### 10.1.1 Summary

For researchers and data scientists who can build and run programs that are three or four pages long, and who want to be more productive and have more confidence in their results, this guide provides a pragmatic, tools-based introduction to program design and maintenance. Unlike books and courses aimed at computer scientists and professional software developers, this guide uses data analysis as a motivating example and assumes that the learner’s ultimate goal is to answer questions rather than ship products.

### 10.1.2 Prerequisites

Learners must be comfortable with the basics of the Unix shell, Python or R, and Git. They will need a personal computer with Internet access, the Bash shell, Python 3, and a GitHub account.

## 10.2 What does “done” look like?

In order to answer the question posed in this section’s title, we need to distinguish between three key ideas. The first is open science, which aims to make



research methods and results available for everyone to read and re-use. The second is reproducible research, which means that anyone with access to the raw materials can easily reproduce the results. Openness and reproducibility are closely related, but are *not* the same thing:

- If you share my data and analysis scripts, but haven't documented the manual steps in the analysis, your work is open but not reproducible.
- If you completely automate the analysis, but it's only available to people in your company or lab, it's reproducible but not open.

The third key idea is sustainability. A piece of software is being sustained if people are using it, fixing it, and improving it rather than replacing it. Sustainability isn't just a property of the software: it also depends on the culture of its actual and potential users. If "share, mend, and extend" is woven into the fabric of their culture, even Fortran-77 can thrive (though of course good tooling and packaging can lower costs and barriers to entry). Conversely, it doesn't matter whether a library has automated tests and is properly packaged if potential users suffer from Not Invented Here syndrome. More importantly, if the software is being maintained by a couple of post-docs who are being paid a fraction of what they could earn in industry, and who have no realistic hope of promotion because their field looks down on tool building, those people will eventually move on and their software will start to suffer from bit rot.

What ties these three ideas together is the notion of computational competence, which is the programming equivalent of good laboratory skills. Software is just another kind of lab equipment; just as an archaeologist should know how to prepare and catalog an artefact, any researcher writing software should know how to make their work reproducible and share it with the world without staying up until dawn.

### Why "Computational Competence"?

The term computational thinking has been widely used since Wing (2006) introduced it a decade ago. It has also been used in such a wide variety of ways that no one really knows what it means. We therefore prefer to talk about computational competence—about someone's ability to do computing well.

## 10.3 What will this course accomplish?

The goal of this course is to help you produce more correct results in less time and with less effort: stakeholders will be confident that you did things the right way, which in turn will allow them to be confident in your results, and you and others will be able to re-use your data, software, and reports instead of constantly rewriting them. To achieve this, we will cover:

- Writing code that is readable, testable, and maintainable

- Automating analyses with build tools
- Checking and demonstrating correctness via automated tests
- Publishing science in the 21st Century
- Using a branch-per-feature workflow, rebasing, and tags to manage work
- Organizing the code, data, results, and reports in a small or medium-sized project

These lessons can be used for self-study by people who are taking part in something like the Insight Data Science Fellows Program, or as part of a one-semester course for graduate students or senior undergraduates. You will know you’re done when:

1. You are reasonably confident that your results are correct. This is not the same as “absolutely sure”: our goal is to make digital work as trustworthy as lab experiments or careful manual analysis.
2. Your software can be used by other people without heroic effort, I.e., people you have never met can find it and figure out how to install it and use it in less time than it would take them to write something themselves.
3. Small changes and extensions are easy so that your software can grow as your problems and questions evolve.

## 10.4 What will we use as running examples?

In order to make this material as accessible as possible, we will use two text processing problems as running examples. The first is an exploration of Zipf’s Law, which states that frequency of a word is inversely proportional to rank, i.e., that the second most common word in some text occurs half as often as most common, the third most common occurs a third as often, and so on. We will write some simple software to test a corpus of text against this rule. The files we will use are taken from the Project Gutenberg and contain this many words:

| Book                            | Words  |
|---------------------------------|--------|
| anne_of_green_gables.txt        | 105642 |
| common_sense.txt                | 24999  |
| count_of_monte_cristo.txt       | 464226 |
| dracula.txt                     | 164424 |
| emma.txt                        | 160458 |
| ethan_frome.txt                 | 37732  |
| frankenstein.txt                | 78098  |
| jane_eyre.txt                   | 188455 |
| life_of_frederick_douglass.txt  | 43789  |
| moby_dick.txt                   | 215830 |
| mysterious_affair_at_styles.txt | 59604  |
| pride_and_prejudice.txt         | 124974 |

| Book                      | Words  |
|---------------------------|--------|
| sense_and_sensibility.txt | 121590 |
| sherlock_holmes.txt       | 107533 |
| time_machine.txt          | 35524  |
| treasure_island.txt       | 71616  |

This is how often the most common words appear in this corpus as a whole:

| Word | Count |
|------|-------|
| the  | 97278 |
| and  | 59385 |
| to   | 56028 |
| of   | 55190 |
| I    | 45680 |
| a    | 40483 |
| in   | 30030 |
| was  | 24512 |
| that | 24386 |
| you  | 22123 |
| it   | 21420 |

The frequencies aren’t an exact match—we would expect about 48,600 occurrences of “and”, for example—but there certainly seems to be a decay curve of some kind. We’ll look more closely at this data as we go along.

The second project is a simple form of computational stylometry. Different writers have different styles; can a computer detect those differences, and if so, can it determine who the likely author of a text actually was? Computational stylometry has been used to explore which parts of Shakespeare’s plays might have been written by other people, which presidential tweets were composed by other people, and who wrote incriminating emails in several high-profile legal cases.

The authors of our books are listed below. Three of them were purportedly written by Jane Austen; we will see if the similarity measures we develop show that.

| Author           | Book                            |
|------------------|---------------------------------|
| Jane Austen      | emma.txt                        |
| Jane Austen      | pride_and_prejudice.txt         |
| Jane Austen      | sense_and_sensibility.txt       |
| Charlotte Brontë | jane_eyre.txt                   |
| Agatha Christie  | mysterious_affair_at_styles.txt |

| Author                      | Book                           |
|-----------------------------|--------------------------------|
| Frederick Douglass          | life_of_frederick_douglass.txt |
| Arthur Conan Doyle          | sherlock_holmes.txt            |
| Alexandre Dumas             | count_of_monte_cristo.txt      |
| Herman Melville             | moby_dick.txt                  |
| Lucy Maud Montgomery        | anne_of_green_gables.txt       |
| Thomas Paine                | common_sense.txt               |
| Mary Wollstonecraft Shelley | frankenstein.txt               |
| Robert Louis Stevenson      | treasure_island.txt            |
| Bram Stoker                 | dracula.txt                    |
| H. G. Wells                 | time_machine.txt               |
| Edith Wharton               | ethan_frome.txt                |

# Chapter 11

## Bash Shell

### 11.1 Questions

- What is a command shell and why would I use one?
- How can I move around on my computer?
- How can I see what files and directories I have?
- How can I specify the location of a file or directory on my computer?
- How can I create, copy, and delete files and directories?
- How can I edit files?
- How can I combine existing commands to do new things?
- How can I perform the same actions on many different files?
- How can I save and re-use commands?
- How can I find files?
- How can I find things in files?

### 11.2 Objectives

- Introduction
  - Explain how the shell relates to the keyboard, the screen, the operating system, and users' programs.
  - Explain when and why command-line interfaces should be used instead of graphical interfaces.
  - Explain the steps in the shell's read-run-print cycle.

- Identify the actual command, options, and filenames in a command-line call.
  - Demonstrate the use of tab completion and explain its advantages.
- Navigating files and directories
  - Explain the similarities and differences between a file and a directory.
  - Translate an absolute path into a relative path and vice versa.
  - Construct absolute and relative paths that identify specific files and directories.
- Working with files and directories
  - Create a directory hierarchy that matches a given diagram.
  - Create files in that hierarchy using an editor or by copying and renaming existing files.
  - Delete, copy and move specified files and/or directories.
- Pipes and filters
  - Redirect a command's output to a file.
  - Process a file instead of keyboard input using redirection.
  - Construct command pipelines with two or more stages.
  - Explain what usually happens if a program or pipeline isn't given any input to process.
  - Explain Unix's 'small pieces, loosely joined' philosophy.
- Loops
  - Write a loop that applies one or more commands separately to each file in a set of files.
  - Trace the values taken on by a loop variable during execution of the loop.
  - Explain the difference between a variable's name and its value.
  - Explain why spaces and some punctuation characters shouldn't be used in file names.
  - Demonstrate how to see what commands have recently been executed.
  - Re-run recently executed commands without retyping them.
- Shell scripts
  - Write a shell script that runs a command or series of commands for a fixed set of files.
  - Run a shell script from the command line.
  - Write a shell script that operates on a set of files defined by the user on the command line.
  - Create pipelines that include shell scripts you, and others, have written.
- Finding things
  - Use **grep** to select lines from text files that match simple patterns.
  - Use **find** to find files whose names match simple patterns.
  - Use the output of one command as the command-line argument(s) to another command.
  - Explain what is meant by 'text' and 'binary' files, and why many common tools don't handle the latter well.

## 11.3 Introduction

At a high level, computers do four things:

- run programs
- store data
- communicate with each other, and
- interact with us

They can do the last of these in many different ways, including through a keyboard and mouse, touch screen interfaces, or using speech recognition systems. While touch and voice interfaces are becoming more commonplace, most interaction is still done using traditional screens, mice, touchpads and keyboards.

The **graphical user interface** (GUI) is the most widely used way to interact with personal computers. We give instructions (to run a program, to copy a file, to create a new folder/directory) with the convenience of a few mouse clicks. This way of interacting with a computer is intuitive and very easy to learn. But this way of giving instructions to a computer scales very poorly if we are to give a large stream of instructions, even if they are similar or identical. For example, suppose we needed to copy the third line of each of a thousand text files stored in thousand different directories and paste it into a single file line by line. Using the traditional GUI approach of mouse clicks, this task would take several hours.

This is where we can take advantage of the shell – a **command-line interface** – to make such repetitive tasks automatic and fast. It can take a single instruction and repeat it (as is or with some modification) as many times as we want. The task in the example above could be accomplished in a few minutes at most.

The heart of a command-line interface is a **read-evaluate-print loop**. It is called so because when you type a command and press Return (also known as Enter) the shell reads your command, evaluates (or “executes”) it, prints the output of your command, then loops back and waits for you to enter another command.

### 11.3.1 The shell

The shell is a program which runs other programs rather than doing calculations itself. Those programs can be as complicated as climate modeling software and as simple as a program that creates a new directory. The simple programs which are used to perform stand alone tasks are usually referred to as commands. The most popular Unix shell is Bash, (the Bourne Again SHell – so-called because it’s derived from a shell written by Stephen Bourne). Bash is the default shell on most modern implementations of Unix and in most packages that provide Unix-like tools for Windows.

When the shell is first opened, you are presented with a **prompt**, indicating that the shell is waiting for input.

```
$
```

The shell typically uses **\$** as the prompt, but may use a different symbol. In the examples for this lesson, we'll show the prompt as **\$**. Most importantly: when typing commands, either from these lessons or from other sources, *do not type the prompt*, only the commands that follow it.

So let's try our first command, which will tell us our username:

```
$ whoami
```

```
beatrice
```

(Beatrice is one of the learner personas described earlier.)

### 11.3.2 Is it difficult?

It is a different model of interacting than a GUI, and that will take some effort – and some time – to learn. A GUI presents you with choices and you select one. With a **command line interface** (CLI) the choices are combinations of commands and parameters, more like words in a language than buttons on a screen. They are not presented to you so you must learn a few, like learning some vocabulary in a new language. But a small number of commands gets you a long way, and we'll cover those essential few in this lesson.

### 11.3.3 Flexibility and automation

The grammar of a shell allows you to combine existing tools into powerful pipelines and handle large volumes of data automatically. Sequences of commands can be written into a *script*, improving the reproducibility of workflows and allowing you to repeat them easily.

In addition, the command line is often the easiest way to interact with remote machines and supercomputers. Familiarity with the shell is near essential to run a variety of specialized tools and resources including high-performance computing systems. As clusters and cloud computing systems become more popular for scientific data crunching, being able to interact with the shell is becoming a necessary skill. We can build on the command-line skills covered here to tackle a wide range of scientific questions and computational challenges.



## 11.4 Navigating Files and Directories

The part of the operating system responsible for managing files and directories is called the **file system**. It organizes our data into files, which hold information, and directories (also called “folders”), which hold files or other directories.

Several commands are frequently used to create, inspect, rename, and delete files and directories. To start exploring them, we’ll go to our open shell window.

First let’s find out where we are by running a command called `pwd` (which stands for “print working directory”). Directories are like *places* – at any time while we are using the shell we are in exactly one place, called our **current working directory**. Commands mostly read and write files in the current working directory, so knowing where you are before running a command is important.

```
$ pwd
```

```
/Users/beatrice
```

Here, the computer’s response is `/Users/beatrice`, which is Beatrice’s **home directory**:

### Home Directory Variation

The home directory path will look different on different operating systems. On Linux it may look like `/home/beatrice`, and on Windows it will be similar to `C:\Documents and Settings\beatrice` or `C:\Users\beatrice`. (Note that it may look slightly different for different versions of Windows.) In future examples, we’ve used Mac output as the default – Linux and Windows output may differ slightly, but should be generally similar.

To understand what a “home directory” is, let’s have a look at how the file system as a whole is organized. For the sake of this example, we’ll be illustrating the filesystem on Beatrice’s computer. After this illustration, you’ll be learning commands to explore your own filesystem, which will be constructed in a similar way, but not be exactly identical.

On Beatrice’s computer, the filesystem looks like this:

*TODO: Figure - The File System ../fig/filesystem.svg*

At the top is the **root directory** that holds everything else. We refer to it using a slash character, `/`, on its own; this is the leading slash in `/Users/beatrice`.

Inside that directory are several other directories: `bin` (which is where some built-in programs are stored), `data` (for miscellaneous data files), `Users` (where users’ personal directories are located), `tmp` (for temporary files that don’t need to be stored long-term), and so on.

We know that our current working directory `/Users/beatrice` is stored inside `/Users` because `/Users` is the first part of its name. Similarly, we know that

`/Users` is stored inside the root directory `/` because its name begins with `/`.

### Slashes

Notice that there are two meanings for the `/` character. When it appears at the front of a file or directory name, it refers to the root directory. When it appears *inside* a name, it's just a separator.

Underneath `/Users`, we find one directory for each user with an account on this machine, namely Beatrice, Jun and Sami.

*TODO: Figure - Home Directories ../fig/home-directories.svg*

The user Jun's files are stored in `/Users/jun`, user Sami's in `/Users/sami`, and Beatrice's in `/Users/beatrice`. Because Beatrice is the user in our examples here, this is why we get `/Users/beatrice` as our home directory. Typically, when you open a new command prompt you will be in your home directory to start.

Now let's learn the command that will let us see the contents of our own filesystem. We can see what's in our home directory by running `ls`, which stands for "listing":

```
$ ls
```

```
Applications Documents Library Music Public
Desktop Downloads Movies Pictures
```

(Again, your results may be slightly different depending on your operating system and how you have customized your filesystem.)

`ls` prints the names of the files and directories in the current directory. We can make its output more comprehensible by using the **-F option** (also known as a **switch** or a **flag**), which tells `ls` to classify the output by adding a marker to file and directory names to indicate what they are: `-` a trailing `/` indicates that this is a directory - `*` indicates an executable (i.e. a program we can run)

Depending on your default options, the shell might also use colors to indicate whether each entry is a file or directory.

```
$ ls -F
```

```
Applications/ Documents/ Library/ Music/ Public/
Desktop/ Downloads/ Movies/ Pictures/
```

Here, we can see that our home directory contains mostly **sub-directories**. Any names in your output that don't have a classification symbol are plain old **files**.

### 11.4.1 General syntax of a shell command

Consider the command below as a general example of a command, which we will dissect into its component parts:

```
$ ls -F /
```

`ls` is the **command**, with an **option** `-F` and an **argument** `/`. We've already encountered options (also called **switches** or **flags**) which either start with a single dash (`-`) or two dashes (`--`), and they change the behaviour of a command. Arguments tell the command what to operate on (e.g. files and directories). Sometimes options and arguments are referred to as **parameters**. A command can be called with more than one option and more than one argument: but a command doesn't always require an argument or an option.

Each part is separated by spaces: if you omit the space between `ls` and `-F` the shell will look for a command called `ls-F`, which doesn't exist. Also, capitalization can be important: `ls -r` is different to `ls -R`.

Putting all that together, our command above gives us a listing of files and directories in the root directory `/`. An example of the output you might get from the above command is given below:

```
$ ls -F /
Applications/      System/
Library/           Users/
Network/           Volumes/
```

### 11.4.2 Getting help

`ls` has lots of other **options**. There are two common ways to find out how to use a command and what options it accepts:

1. We can pass a `--help` option to the command, such as: `$ ls --help`
2. We can read its manual with `man`, such as: `$ man ls`

Depending on your environment you might find that only one of these works (either `man` or `--help`). We'll describe both ways below.

#### 11.4.2.1 The `--help` option

Many bash commands, and programs that people have written that can be run from within bash, support a `--help` option to display more information on how to use the command or program.

```
$ ls --help
```

Usage: `ls [OPTION]... [FILE]...`

List information about the FILES (the current directory by default).

Sort entries alphabetically if none of `-cftuvSUX` nor `--sort` is specified.

Mandatory arguments to long options are mandatory for short options too.

|  |  |
|--|--|
| <code>-a, --all</code>                                 | do not ignore entries starting with <code>.</code>   |
| <code>-A, --almost-all</code>                          | do not list implied <code>.</code> and <code>..</code>   |
| <code>--author</code>                                  | with <code>-l</code> , print the author of each file   |
| <code>-b, --escape</code>                              | print C-style escapes for nongraphic characters  |
| <code>--block-size=SIZE</code>                         | scale sizes by SIZE before printing them; e.g.,<br>' <code>--block-size=M</code> ' prints sizes in units of<br>1,048,576 bytes; see SIZE format below  |
| <code>-B, --ignore-backups</code>                      | do not list implied entries ending with <code>~</code>   |
| <code>-c</code>  | with <code>-lt</code> : sort by, and show, ctime (time of last<br>modification of file status information);<br>with <code>-l</code> : show ctime and sort by name;<br>otherwise: sort by ctime, newest first |
| <code>-C</code>  | list entries by columns  |
| <code>--color[=WHEN]</code>                            | colorize the output; WHEN can be 'always' (default<br>if omitted), 'auto', or 'never'; more info below   |
| <code>-d, --directory</code>                           | list directories themselves, not their contents  |
| <code>-D, --dired</code>                               | generate output designed for Emacs' dired mode   |
| <code>-f</code>  | do not sort, enable <code>-aU</code> , disable <code>-ls --color</code>  |
| <code>-F, --classify</code>                            | append indicator (one of <code>*/=&gt;@ </code> ) to entries   |
| <code>--file-type</code>                               | likewise, except do not append <code>'*</code>   |
| <code>--format=WORD</code>                             | across <code>-x</code> , commas <code>-m</code> , horizontal <code>-x</code> , long <code>-l</code> ,<br>single-column <code>-1</code> , verbose <code>-l</code> , vertical <code>-C</code>                  |
| <code>--full-time</code>                               | like <code>-l --time-style=full-iso</code>   |
| <code>-g</code>  | like <code>-l</code> , but do not list owner   |
| <code>--group-directories-first</code>                 | group directories before files;<br>can be augmented with a <code>--sort</code> option, but any<br>use of <code>--sort=none (-U)</code> disables grouping   |
| <code>-G, --no-group</code>                            | in a long listing, don't print group names   |
| <code>-h, --human-readable</code>                      | with <code>-l</code> and/or <code>-s</code> , print human readable sizes<br>(e.g., 1K 234M 2G)   |
| <code>--si</code>                                      | likewise, but use powers of 1000 not 1024  |
| <code>-H, --dereference-command-line</code>            | follow symbolic links listed on the command line   |
| <code>--dereference-command-line-symlink-to-dir</code> | follow each command line symbolic link<br>that points to a directory   |
| <code>--hide=PATTERN</code>                            | do not list implied entries matching shell PATTERN<br>(overridden by <code>-a</code> or <code>-A</code> )  |
| <code>--indicator-style=WORD</code>                    | append indicator with style WORD to entry names:<br>none (default), slash ( <code>-p</code> ),   |

|                             |  |
|-----------------------------|--|
|                             | file-type (--file-type), classify (-F)   |
| -i, --inode                 | print the index number of each file  |
| -I, --ignore=PATTERN        | do not list implied entries matching shell PATTERN   |
| -k, --kibibytes             | default to 1024-byte blocks for disk usage   |
| -l                          | use a long listing format  |
| -L, --dereference           | when showing file information for a symbolic link, show information for the file the link references rather than for the link itself   |
| -m                          | fill width with a comma separated list of entries  |
| -n, --numeric-uid-gid       | like -l, but list numeric user and group IDs   |
| -N, --literal               | print raw entry names (don't treat e.g. control characters specially)  |
| -o                          | like -l, but do not list group information   |
| -p, --indicator-style=slash | append / indicator to directories  |
| -q, --hide-control-chars    | print ? instead of nongraphic characters   |
| --show-control-chars        | show nongraphic characters as-is (the default, unless program is 'ls' and output is a terminal)  |
| -Q, --quote-name            | enclose entry names in double quotes   |
| --quoting-style=WORD        | use quoting style WORD for entry names:<br>literal, locale, shell, shell-always,<br>shell-escape, shell-escape-always, c, escape   |
| -r, --reverse               | reverse order while sorting  |
| -R, --recursive             | list subdirectories recursively  |
| -s, --size                  | print the allocated size of each file, in blocks   |
| -S                          | sort by file size, largest first   |
| --sort=WORD                 | sort by WORD instead of name: none (-U), size (-S),<br>time (-t), version (-v), extension (-X)   |
| --time=WORD                 | with -l, show time as WORD instead of default<br>modification time: atime or access or use (-u);<br>ctime or status (-c); also use specified time<br>as sort key if --sort=time (newest first)   |
| --time-style=STYLE          | with -l, show times using style STYLE:<br>full-iso, long-iso, iso, locale, or +FORMAT;<br>FORMAT is interpreted like in 'date'; if FORMAT<br>is FORMAT1<newline>FORMAT2, then FORMAT1 applies<br>to non-recent files and FORMAT2 to recent files;<br>if STYLE is prefixed with 'posix-', STYLE<br>takes effect only outside the POSIX locale |
| -t                          | sort by modification time, newest first  |
| -T, --tabsize=COLS          | assume tab stops at each COLS instead of 8   |
| -u                          | with -lt: sort by, and show, access time;<br>with -l: show access time and sort by name;<br>otherwise: sort by access time, newest first   |
| -U                          | do not sort; list entries in directory order   |
| -v                          | natural sort of (version) numbers within text  |

```

-w, --width=COLS      set output width to COLS. 0 means no limit
-x                    list entries by lines instead of by columns
-X                    sort alphabetically by entry extension
-Z, --context          print any security context of each file
-1                    list one file per line. Avoid '\n' with -q or -b
    --help            display this help and exit
    --version          output version information and exit

```

The `SIZE` argument is an integer and optional unit (example: 10K is 10\*1024). Units are K,M,G,T,P,E,Z,Y (powers of 1024) or KB,MB,... (powers of 1000).

Using color to distinguish file types is disabled both by default and with `--color=never`. With `--color=auto`, `ls` emits color codes only when standard output is connected to a terminal. The `LS_COLORS` environment variable can change the settings. Use the `dircolors` command to set it.

Exit status:

```

0  if OK,
1  if minor problems (e.g., cannot access subdirectory),
2  if serious trouble (e.g., cannot access command-line argument).

```

GNU coreutils online help: <<http://www.gnu.org/software/coreutils/>>  
 Full documentation at: <<http://www.gnu.org/software/coreutils/ls>>  
 or available locally via: `info '(coreutils) ls invocation'`

### Unsupported command-line options

If you try to use an option that is not supported, `ls` and other commands will usually print an error message similar to:

```

$ ls -j

ls: invalid option -- 'j'
Try 'ls --help' for more information.

```

#### 11.4.2.2 The `man` command

The other way to learn about `ls` is to type

```
$ man ls
```

This will display a description of the `ls` command and its options and, if you're lucky, some examples of how to use it.

To navigate through the `man` pages, you may use `↑` and `↓` to move line-by-line, or try `B` and `Spacebar` to skip up and down by a full page. To search for a character or word in the `man` pages, use `/` followed by the character or word you are searching for. Sometimes a search will result in multiple hits. If so, you

can move between hits using N (for moving forward) and Shift+N (for moving backward).

To quit the `man` pages, press Q.

#### Assistance on the web

Of course there is a third way to access help for commands: searching the internet via your web browser. In many cases the first results from your search will be Stack Overflow pages where someone has already asked (and hopefully had answered) a question similar to yours.

GNU provides links to its manuals including the core GNU utilities, which covers many commands introduced within this lesson. There's also a community effort known as the TLDR pages that aims to simplify the default man pages with practical examples.

### 11.4.3 Looking around and moving about

We can also use `ls` to see the contents of a different directory. Let's take a look at our `Desktop` directory by running `ls -F Desktop`, i.e., the command `ls` with the `-F` option and the argument `Desktop`. The argument `Desktop` tells `ls` that we want a listing of something other than our current working directory:

```
$ ls -F Desktop
```

```
data-shell/
```

Your output should be a list of all the files and sub-directories on your Desktop, including the `data-shell` directory you downloaded at the setup for this lesson. Take a look at your Desktop to confirm that your output is accurate.

Now that we know the `data-shell` directory is located on our Desktop, we can do two things.

First, we can look at its contents, using the same strategy as before, passing a directory name to `ls`:

```
$ ls -F Desktop/data-shell
```

```
creatures/      molecules/      notes.txt      solar.pdf
data/           north-pacific-gyre/  pizza.cfg      writing/
```

Second, we can actually change our location to a different directory, so we are no longer located in our home directory.

The command to change locations is `cd` followed by a directory name to change our working directory. `cd` stands for “change directory”, which is a bit misleading: the command doesn't change the directory, it changes the shell's idea of what directory we are in.

Let's say we want to move to the `data` directory we saw above. We can use the following series of commands to get there:

```
$ cd Desktop
$ cd data-shell
$ cd data
```

These commands will move us from our home directory onto our Desktop, then into the `data-shell` directory, then into the `data` directory. You will notice that `cd` doesn't print anything. This is normal. Many shell commands will not output anything to the screen when successfully executed. But if we run `pwd` after it, we can see that we are now in `/Users/beatrice/Desktop/data-shell/data`. If we run `ls` without arguments now, it lists the contents of `/Users/beatrice/Desktop/data-shell/data`, because that's where we now are:

```
$ pwd
/Users/beatrice/Desktop/data-shell/data
$ ls -F
amino-acids.txt  elements/      pdb/           salmon.txt
animals.txt      morse.txt      planets.txt    sunspot.txt
```

We now know how to go down the directory tree, but how do we go up? We might try the following:

```
$ cd data-shell
-bash: cd: data-shell: No such file or directory
```

But we get an error! Why is this?

With our methods so far, `cd` can only see sub-directories inside your current directory. There are different ways to see directories above your current location; we'll start with the simplest.

There is a shortcut in the shell to move up one directory level that looks like this:

```
$ cd ..
```

`..` is a special directory name meaning "the directory containing this one", or more succinctly, the **parent** of the current directory. Sure enough, if we run `pwd` after running `cd ..`, we're back in `/Users/beatrice/Desktop/data-shell`:

```
$ pwd
/Users/beatrice/Desktop/data-shell
```

The special directory `..` doesn't usually show up when we run `ls`. If we want to display it, we can give `ls` the `-a` option:

```
$ ls -F -a
```



```
./  .bash_profile  data/      north-pacific-gyre/  pizza.cfg  thesis/
../  creatures/    molecules/  notes.txt           solar.pdf  writing/
```

`-a` stands for “show all”; it forces `ls` to show us file and directory names that begin with `.`, such as `..` (which, if we’re in `/Users/beatrice`, refers to the `/Users` directory) As you can see, it also displays another special directory that’s just called `.`, which means “the current working directory”. It may seem redundant to have a name for it, but we’ll see some uses for it soon.

Note that in most command line tools, multiple options can be combined with a single `-` and no spaces between the options: `ls -F -a` is equivalent to `ls -Fa`.

### Other Hidden Files

In addition to the hidden directories `..` and `.`, you may also see a file called `.bash_profile`. This file usually contains shell configuration settings. You may also see other files and directories beginning with `..`. These are usually files and directories that are used to configure different programs on your computer. The prefix `.` is used to prevent these configuration files from cluttering the terminal when a standard `ls` command is used.

### Orthogonality

The special names `.` and `..` don’t belong to `cd`; they are interpreted the same way by every program. For example, if we are in `/Users/beatrice/data`, the command `ls ..` will give us a listing of `/Users/beatrice`. When the meanings of the parts are the same no matter how they’re combined, programmers say they are **orthogonal**: Orthogonal systems tend to be easier for people to learn because there are fewer special cases and exceptions to keep track of.

So these are the basic commands for navigating the filesystem on your computer: `pwd`, `ls` and `cd`. Let’s explore some variations on those commands. What happens if you type `cd` on its own, without giving a directory?

```
$ cd
```

How can you check what happened? `pwd` gives us the answer:

```
$ pwd
```

```
/Users/beatrice
```

It turns out that `cd` without an argument will return you to your home directory, which is great if you’ve gotten lost in your own filesystem.

Let’s try returning to the `data` directory from before. Last time, we used three commands, but we can actually string together the list of directories to move to `data` in one step:

```
$ cd Desktop/data-shell/data
```

Check that we've moved to the right place by running `pwd` and `ls -F`

If we want to move up one level from the `data` directory, we could use `cd ..`. But there is another way to move to any directory, regardless of your current location.

So far, when specifying directory names, or even a directory path (as above), we have been using **relative paths**. When you use a relative path with a command like `ls` or `cd`, it tries to find that location from where we are, rather than from the root of the file system.

However, it is possible to specify the **absolute path** to a directory by including its entire path from the root directory, which is indicated by a leading slash. The leading `/` tells the computer to follow the path from the root of the file system, so it always refers to exactly one directory, no matter where we are when we run the command.

This allows us to move to our `data-shell` directory from anywhere on the filesystem (including from inside `data`). To find the absolute path we're looking for, we can use `pwd` and then extract the piece we need to move to `data-shell`.

```
$ pwd
```

```
/Users/beatrice/Desktop/data-shell/data
```

```
$ cd /Users/beatrice/Desktop/data-shell
```

Run `pwd` and `ls -F` to ensure that we're in the directory we expect.

### Two More Shortcuts

The shell interprets the character `~` (tilde) at the start of a path to mean “the current user's home directory”. For example, if Beatrice's home directory is `/Users/beatrice`, then `~/data` is equivalent to `/Users/beatrice/data`. This only works if it is the first character in the path: `here/there/~/elsewhere` is *not* `here/there/Users/beatrice/elsewhere`.

Another shortcut is the `-` (dash) character. `cd` will translate `-` into *the previous directory you were in*, which is faster than having to remember, then type, the full path. This is a *very* efficient way of moving back and forth between directories. The difference between `cd ..` and `cd -` is that the former brings you *up*, while the latter brings you *back*. You can think of it as the *Last Channel* button on a TV remote.

## 11.5 Working With Files and Directories

### 11.5.1 Creating directories

We now know how to explore files and directories, but how do we create them in the first place?

Step one: see where we are and what we already have. Let's go back to our `data-shell` directory on the Desktop and use `ls -F` to see what it contains:

```
$ pwd
/Users/beatrice/Desktop/data-shell
$ ls -F
creatures/ data/ molecules/ north-pacific-gyre/ notes.txt pizza.cfg solar.pdf writing/
```

Let's create a new directory called `thesis` using the command `mkdir thesis` (which has no output):

```
$ mkdir thesis
```

As you might guess from its name, `mkdir` means “make directory”. Since `thesis` is a relative path (i.e., does not have a leading slash, like `/what/ever/thesis`), the new directory is created in the current working directory:

```
$ ls -F
creatures/ data/ molecules/ north-pacific-gyre/ notes.txt pizza.cfg solar.pdf thesis/ wri
```

#### Two ways of doing the same thing

Using the shell to create a directory is no different than using a file explorer. If you open the current directory using your operating system's graphical file explorer, the `thesis` directory will appear there too. While the shell and the file explorer are two different ways of interacting with the files, the files and directories themselves are the same.

#### Good names for files and directories

Complicated names of files and directories can make your life painful when working on the command line. Here we provide a few useful tips for the names of your files.

1. Don't use spaces.

Spaces can make a name more meaningful, but since spaces are used to separate arguments on the command line it is better to avoid them in names of files and directories. You can use `-` or `_` instead (e.g. `north-pacific-gyre/` rather than `north pacific gyre/`).

2. Don't begin the name with - (dash).

Commands treat names starting with - as options.

3. Stick with letters, numbers, . (period or 'full stop'), - (dash) and \_ (underscore).

Many other characters have special meanings on the command line. We will learn about some of these during this lesson. There are special characters that can cause your command to not work as expected and can even result in data loss.

If you need to refer to names of files or directories that have spaces or other special characters, you should surround the name in quotes ("").

Since we've just created the `thesis` directory, there's nothing in it yet:

```
$ ls -F thesis
```

### 11.5.2 Create a text file

Let's change our working directory to `thesis` using `cd`, then run a text editor called Nano to create a file called `draft.txt`:

```
$ cd thesis
$ nano draft.txt
```

#### Which Editor?

When we say, “`nano` is a text editor,” we really do mean “text”: it can only work with plain character data, not tables, images, or any other human-friendly media. We use it in examples because it is one of the least complex text editors. However, because of this trait, it may not be powerful enough or flexible enough for the work you need to do after this workshop. On Unix systems (such as Linux and Mac OS X), many programmers use Emacs or Vim (both of which require more time to learn), or a graphical editor such as Sublime Text. On Windows, a popular editor is Notepad++. Windows also has a built-in editor called `notepad` that can be run from the command line in the same way as `nano` for the purposes of this lesson.

No matter what editor you use, you will need to know where it searches for and saves files. If you start it from the shell, it will (probably) use your current working directory as its default location. If you use your computer's start menu, it may want to save files in your desktop or documents directory instead. You can change this by navigating to another directory the first time you “Save As...”

Let's type in a few lines of text. Once we're happy with our text, we can press **Ctrl+O** (press the **Ctrl** or **Control** key and, while holding it down, press the **O** key) to write our data to disk (we'll be asked what file we want to save this to: press **Return** to accept the suggested default of **draft.txt**).

*TODO: Figure - Nano in Action ../fig/nano-screenshot.png*

Once our file is saved, we can use **Ctrl-X** to quit the editor and return to the shell.

### Control, Ctrl, or ^ Key

The Control key is also called the "Ctrl" key. There are various ways in which using the Control key may be described. For example, you may see an instruction to press the Control key and, while holding it down, press the X key, described as any of:

- **Control-X**
- **Control+X**
- **Ctrl-X**
- **Ctrl+X**
- **^X**
- **C-x**

In nano, along the bottom of the screen you'll see **^G Get Help ^O WriteOut**. This means that you can use **Control-G** to get help and **Control-O** to save your file.

**nano** doesn't leave any output on the screen after it exits, but **ls** now shows that we have created a file called **draft.txt**:

```
$ ls
draft.txt
```

### What's In A Name?

You may have noticed that all of Beatrice's files are named "something dot something", and in this part of the lesson, we always used the extension **.txt**. This is just a convention: we can call a file **mythesis** or almost anything else we want. However, most people use two-part names most of the time to help them (and their programs) tell different kinds of files apart. The second part of such a name is called the **filename extension**, and indicates what type of data the file holds: **.txt** signals a plain text file, **.pdf** indicates a PDF document, **.cfg** is a configuration file full of parameters for some program or other, **.png** is a PNG image, and so on.

This is just a convention, albeit an important one. Files contain bytes: it's up to us and our programs to interpret those bytes according to the rules for plain text files, PDF documents, configuration files, images, and so on.

Naming a PNG image of a whale as `whale.mp3` doesn't somehow magically turn it into a recording of whalesong, though it *might* cause the operating system to try to open it with a music player when someone double-clicks it.

### 11.5.3 Moving files and directories

Returning to the `data-shell` directory,

```
cd ~/Desktop/data-shell/
```

In our `thesis` directory we have a file `draft.txt` which isn't a particularly informative name, so let's change the file's name using `mv`, which is short for "move":

```
$ mv thesis/draft.txt thesis/quotes.txt
```

The first argument tells `mv` what we're "moving", while the second is where it's to go. In this case, we're moving `thesis/draft.txt` to `thesis/quotes.txt`, which has the same effect as renaming the file. Sure enough, `ls` shows us that `thesis` now contains one file called `quotes.txt`:

```
$ ls thesis
```

```
quotes.txt
```

One has to be careful when specifying the target file name, since `mv` will silently overwrite any existing file with the same name, which could lead to data loss. An additional option, `mv -i` (or `mv --interactive`), can be used to make `mv` ask you for confirmation before overwriting.

Note that `mv` also works on directories.

Let's move `quotes.txt` into the current working directory. We use `mv` once again, but this time we'll just use the name of a directory as the second argument to tell `mv` that we want to keep the filename, but put the file somewhere new. (This is why the command is called "move".) In this case, the directory name we use is the special directory name `.` that we mentioned earlier.

```
$ mv thesis/quotes.txt .
```

The effect is to move the file from the directory it was in to the current working directory. `ls` now shows us that `thesis` is empty:

```
$ ls thesis
```

Further, `ls` with a filename or directory name as an argument only lists that file or directory. We can use this to see that `quotes.txt` is still in our current directory:

```
$ ls quotes.txt
```

quotes.txt

### 11.5.4 Copying files and directories

The `cp` command works very much like `mv`, except it copies a file instead of moving it. We can check that it did the right thing using `ls` with two paths as arguments — like most Unix commands, `ls` can be given multiple paths at once:

```
$ cp quotes.txt thesis/quotations.txt
$ ls quotes.txt thesis/quotations.txt

quotes.txt  thesis/quotations.txt
```

We can also copy a directory and all its contents by using the recursive option `-r`, e.g. to back up a directory:

```
$ cp -r thesis thesis_backup
```

We can check the result by listing the contents of both the `thesis` and `thesis_backup` directory:

```
$ ls thesis thesis_backup

thesis:
quotations.txt

thesis_backup:
quotations.txt
```

### 11.5.5 Removing files and directories

Returning to the `data-shell` directory, let's tidy up this directory by removing the `quotes.txt` file we created. The Unix command we'll use for this is `rm` (short for 'remove'):

```
$ rm quotes.txt
```

We can confirm the file has gone using `ls`:

```
$ ls quotes.txt
```

```
ls: cannot access 'quotes.txt': No such file or directory
```

#### Deleting Is Forever

The Unix shell doesn't have a trash bin that we can recover deleted files from (though most graphical interfaces to Unix do). Instead, when we delete files, they are unlinked from the file system so that their storage space on disk can be recycled. Tools for finding and

recovering deleted files do exist, but there's no guarantee they'll work in any particular situation, since the computer may recycle the file's disk space right away.

If we try to remove the `thesis` directory using `rm thesis`, we get an error message:

```
$ rm thesis
```

```
rm: cannot remove `thesis': Is a directory
```

This happens because `rm` by default only works on files, not directories.

`rm` can remove a directory *and all its contents* if we use the recursive option `-r`, and it will do so *without any confirmation prompts*:

```
$ rm -r thesis
```

Given that there is no way to retrieve files deleted using the shell, `rm -r` *should be used with great caution* (you might consider adding the interactive option `rm -r -i`).

### 11.5.6 Operations with multiple files and directories

Oftentimes one needs to copy or move several files at once. This can be done by providing a list of individual filenames, or specifying a naming pattern using wildcards.

#### Wildcards

`*` is a **wildcard**, which matches zero or more characters. Let's consider the `data-shell/molecules` directory: `*.pdb` matches `ethane.pdb`, `propane.pdb`, and every file that ends with `.pdb`. On the other hand, `p*.pdb` only matches `pentane.pdb` and `propane.pdb`, because the `'p'` at the front only matches filenames that begin with the letter `'p'`.

`?` is also a wildcard, but it matches exactly one character. So `?ethane.pdb` would match `methane.pdb` whereas `*ethane.pdb` matches both `ethane.pdb`, and `methane.pdb`.

Wildcards can be used in combination with each other e.g. `???ane.pdb` matches three characters followed by `ane.pdb`, giving `cubane.pdb` `ethane.pdb` `octane.pdb`.

When the shell sees a wildcard, it expands the wildcard to create a list of matching filenames *before* running the command that was asked for. As an exception, if a wildcard expression does not match any file, Bash will pass the expression as an argument to the command as it is. For example typing `ls *.pdf` in the `molecules` directory (which contains only files with names ending with `.pdb`) results



in an error message that there is no file called `*.pdf`. However, generally commands like `wc` and `ls` see the lists of file names matching these expressions, but not the wildcards themselves. It is the shell, not the other programs, that deals with expanding wildcards, and this is another example of orthogonal design.

## 11.6 Pipes and Filters

Now that we know a few basic commands, we can finally look at the shell's most powerful feature: the ease with which it lets us combine existing programs in new ways. We'll start with a directory called `molecules` that contains six files describing some simple organic molecules. The `.pdb` extension indicates that these files are in Protein Data Bank format, a simple text format that specifies the type and position of each atom in the molecule.

```
$ ls molecules
```

```
cubane.pdb    ethane.pdb    methane.pdb
octane.pdb    pentane.pdb    propane.pdb
```

Let's go into that directory with `cd` and run the command `wc *.pdb`. `wc` is the "word count" command: it counts the number of lines, words, and characters in files (from left to right, in that order).

The `*` in `*.pdb` matches zero or more characters, so the shell turns `*.pdb` into a list of all `.pdb` files in the current directory:

```
$ cd molecules
```

```
$ wc *.pdb
```

```
20 156 1158 cubane.pdb
12 84 622 ethane.pdb
9 57 422 methane.pdb
30 246 1828 octane.pdb
21 165 1226 pentane.pdb
15 111 825 propane.pdb
107 819 6081 total
```

If we run `wc -l` instead of just `wc`, the output shows only the number of lines per file:

```
$ wc -l *.pdb
```

```
20 cubane.pdb
12 ethane.pdb
9 methane.pdb
30 octane.pdb
21 pentane.pdb
```

```
15  propane.pdb
107 total
```

### Why Isn't It Doing Anything?

What happens if a command is supposed to process a file, but we don't give it a filename? For example, what if we type:

```
$ wc -l
```

but don't type `*.pdb` (or anything else) after the command? Since it doesn't have any filenames, `wc` assumes it is supposed to process standard input, so it just sits there and waits for us to give it some data interactively. From the outside, though, all we see is it sitting there: the command doesn't appear to do anything.

If you make this kind of mistake, you can escape out of this state by holding down the control key (Ctrl) and typing the letter C once and letting go of the Ctrl key. Ctrl+C

We can also use `-w` to get only the number of words, or `-c` to get only the number of characters.

Which of these files contains the fewest lines? It's an easy question to answer when there are only six files, but what if there were 6000? Our first step toward a solution is to run the command:

```
$ wc -l *.pdb > lengths.txt
```

The greater than symbol, `>`, tells the shell to **redirect** the command's output to a file instead of printing it to the screen. (This is why there is no screen output: everything that `wc` would have printed has gone into the file `lengths.txt` instead.) The shell will create the file if it doesn't exist. If the file exists, it will be silently overwritten, which may lead to data loss and thus requires some caution. `ls lengths.txt` confirms that the file exists:

```
$ ls lengths.txt
lengths.txt
```

We can now send the content of `lengths.txt` to the screen using `cat lengths.txt`. `cat` stands for "concatenate": it prints the contents of files one after another. There's only one file in this case, so `cat` just shows us what it contains:

```
$ cat lengths.txt
20  cubane.pdb
12  ethane.pdb
9   methane.pdb
30  octane.pdb
21  pentane.pdb
15  propane.pdb
```

```
107 total
```

### Output Page by Page

We'll continue to use `cat` in this lesson, for convenience and consistency, but it has the disadvantage that it always dumps the whole file onto your screen. More useful in practice is the command `less`, which you use with `less lengths.txt`. This displays a screenful of the file, and then stops. You can go forward one screenful by pressing the spacebar, or back one by pressing `b`. Press `q` to quit.

Now let's use the `sort` command to sort its contents.

We will also use the `-n` option to specify that the sort is numerical instead of alphanumerical. This does *not* change the file; instead, it sends the sorted result to the screen:

```
$ sort -n lengths.txt
```

```
9  methane.pdb
12 ethane.pdb
15 propane.pdb
20 cubane.pdb
21 pentane.pdb
30 octane.pdb
107 total
```

We can put the sorted list of lines in another temporary file called `sorted-lengths.txt` by putting `> sorted-lengths.txt` after the command, just as we used `> lengths.txt` to put the output of `wc` into `lengths.txt`. Once we've done that, we can run another command called `head` to get the first few lines in `sorted-lengths.txt`:

```
$ sort -n lengths.txt > sorted-lengths.txt
$ head -n 1 sorted-lengths.txt
```

```
9  methane.pdb
```

Using `-n 1` with `head` tells it that we only want the first line of the file; `-n 20` would get the first 20, and so on. Since `sorted-lengths.txt` contains the lengths of our files ordered from least to greatest, the output of `head` must be the file with the fewest lines.

### Redirecting to the same file

It's a very bad idea to try redirecting the output of a command that operates on a file to the same file. For example:

```
$ sort -n lengths.txt > lengths.txt
```

Doing something like this may give you incorrect results and/or delete the contents of `lengths.txt`.

If you think this is confusing, you're in good company: even once you understand what `wc`, `sort`, and `head` do, all those intermediate files make it hard to follow what's going on. We can make it easier to understand by running `sort` and `head` together:

```
$ sort -n lengths.txt | head -n 1

9  methane.pdb
```

The vertical bar, `|`, between the two commands is called a **pipe**. It tells the shell that we want to use the output of the command on the left as the input to the command on the right. The computer might create a temporary file if it needs to, or copy data from one program to the other in memory, or something else entirely; we don't have to know or care.

Nothing prevents us from chaining pipes consecutively. That is, we can for example send the output of `wc` directly to `sort`, and then the resulting output to `head`. Thus we first use a pipe to send the output of `wc` to `sort`:

```
$ wc -l *.pdb | sort -n

9  methane.pdb
12 ethane.pdb
15 propane.pdb
20 cubane.pdb
21 pentane.pdb
30 octane.pdb
107 total
```

And now we send the output of this pipe, through another pipe, to `head`, so that the full pipeline becomes:

```
$ wc -l *.pdb | sort -n | head -n 1

9  methane.pdb
```

This is exactly like a mathematician nesting functions like  $\log(3x)$  and saying “the log of three times  $x$ ”. In our case, the calculation is “head of sort of line count of `*.pdb`”.

Here's what actually happens behind the scenes when we create a pipe. When a computer runs a program — any program — it creates a **process** in memory to hold the program's software and its current state. Every process has an input channel called **standard input**. (By this point, you may be surprised that the name is so memorable, but don't worry: most Unix programmers call it “stdin”). Every process also has a default output channel called **standard output** (or “stdout”). A second output channel called **standard error** (stderr) also exists. This channel is typically used for error or diagnostic messages, and it allows a user to pipe the output of one program into another while still receiving error messages in the terminal.

The shell is actually just another program. Under normal circumstances, whatever we type on the keyboard is sent to the shell on its standard input, and whatever it produces on standard output is displayed on our screen. When we tell the shell to run a program, it creates a new process and temporarily sends whatever we type on our keyboard to that process's standard input, and whatever the process sends to standard output to the screen.

Here's what happens when we run `wc -l *.pdb > lengths.txt`. The shell starts by telling the computer to create a new process to run the `wc` program. Since we've provided some filenames as arguments, `wc` reads from them instead of from standard input. And since we've used `>` to redirect output to a file, the shell connects the process's standard output to that file.

If we run `wc -l *.pdb | sort -n` instead, the shell creates two processes (one for each process in the pipe) so that `wc` and `sort` run simultaneously. The standard output of `wc` is fed directly to the standard input of `sort`; since there's no redirection with `>`, `sort`'s output goes to the screen. And if we run `wc -l *.pdb | sort -n | head -n 1`, we get three processes with data flowing from the files, through `wc` to `sort`, and from `sort` through `head` to the screen.

*TODO: Figure - Redirects and Pipes ../fig/redirects-and-pipes.png*

This simple idea is why Unix has been so successful. Instead of creating enormous programs that try to do many different things, Unix programmers focus on creating lots of simple tools that each do one job well, and that work well with each other. This programming model is called "pipes and filters". We've already seen pipes; a **filter** is a program like `wc` or `sort` that transforms a stream of input into a stream of output. Almost all of the standard Unix tools can work this way: unless told to do otherwise, they read from standard input, do something with what they've read, and write to standard output.

The key is that any program that reads lines of text from standard input and writes lines of text to standard output can be combined with every other program that behaves this way as well. You can *and should* write your programs this way so that you and other people can put those programs into pipes to multiply their power.

## 11.7 Loops

**Loops** are a programming construct which allow us to repeat a command or set of commands for each item in a list. As such they are key to productivity improvements through automation. Similar to wildcards and tab completion, using loops also reduces the amount of typing required (and hence reduces the number of typing mistakes).

Suppose we have several hundred genome data files named `basilisk.dat`, `unicorn.dat`, and so on. For this example, we'll use the `creatures` directory

which only has two example files, but the principles can be applied to many many more files at once. We would like to print out the classification for each species, which is given on the second line of the file. For each file, we would need to execute the command `head -n 2` and pipe this to `tail -n 1`. We'll use a loop to solve this problem, but first let's look at the general form of a loop:

```
for thing in list_of_things
do
    operation_using $thing    # Indentation within the loop is not required, but aids
done
```

and we can apply this to our example like this:

```
$ for filename in basilisk.dat unicorn.dat
> do
>   head -n 2 $filename | tail -n 1
> done
```

```
CLASSIFICATION: basiliscus vulgaris
```

```
CLASSIFICATION: equus monoceros
```

### Follow the Prompt

The shell prompt changes from `$` to `>` and back again as we were typing in our loop. The second prompt, `>`, is different to remind us that we haven't finished typing a complete command yet. A semicolon, `;`, can be used to separate two commands written on a single line.

When the shell sees the keyword `for`, it knows to repeat a command (or group of commands) once for each item in a list. Each time the loop runs (called an iteration), an item in the list is assigned in sequence to the **variable**, and the commands inside the loop are executed, before moving on to the next item in the list. Inside the loop, we call for the variable's value by putting `$` in front of it. The `$` tells the shell interpreter to treat the variable as a variable name and substitute its value in its place, rather than treat it as text or an external command.

In this example, the list is two filenames: `basilisk.dat` and `unicorn.dat`. Each time the loop iterates, it will assign a file name to the variable `filename` and run the `head` command. The first time through the loop, `$filename` is `basilisk.dat`. The interpreter runs the command `head` on `basilisk.dat`, and then prints the first three lines of `basilisk.dat`. For the second iteration, `$filename` becomes `unicorn.dat`. This time, the shell runs `head` on `unicorn.dat` and prints the first three lines of `unicorn.dat`. Since the list was only two items, the shell exits the `for` loop.

### Same Symbols, Different Meanings

Here we see `>` being used as a shell prompt, whereas `>` is also used to redirect output. Similarly, `$` is used as a shell prompt, but, as we saw earlier, it is also used to ask the shell to get the value of a variable.

If the *shell* prints `>` or `$` then it expects you to type something, and the symbol is a prompt.

If *you* type `>` or `$` yourself, it is an instruction from you that the shell should redirect output or get the value of a variable.

When using variables it is also possible to put the names into curly braces to clearly delimit the variable name: `$filename` is equivalent to `${filename}`, but is different from `${file}name`. You may find this notation in other people's programs.

We have called the variable in this loop `filename` in order to make its purpose clearer to human readers. The shell itself doesn't care what the variable is called; if we wrote this loop as:

```
$ for x in basilisk.dat unicorn.dat
> do
>   head -n 2 $x | tail -n 1
> done
```

or:

```
$ for temperature in basilisk.dat unicorn.dat
> do
>   head -n 2 $temperature | tail -n 1
> done
```

it would work exactly the same way. *Don't do this.* Programs are only useful if people can understand them, so meaningless names (like `x`) or misleading names (like `temperature`) increase the odds that the program won't do what its readers think it does.

Let's continue with our example in the `data-shell/creatures` directory. Here's a slightly more complicated loop:

```
$ for filename in *.dat
> do
>   echo $filename
>   head -n 100 $filename | tail -n 20
> done
```

The shell starts by expanding `*.dat` to create the list of files it will process. The **loop body** then executes two commands for each of those files. The first, `echo`, just prints its command-line arguments to standard output. For example:

```
$ echo hello there
```

prints:

hello there

In this case, since the shell expands `$filename` to be the name of a file, `echo $filename` just prints the name of the file. Note that we can't write this as:

```
$ for filename in *.dat
> do
>     $filename
>     head -n 100 $filename | tail -n 20
> done
```

because then the first time through the loop, when `$filename` expanded to `basilisk.dat`, the shell would try to run `basilisk.dat` as a program. Finally, the `head` and `tail` combination selects lines 81-100 from whatever file is being processed (assuming the file has at least 100 lines).

### Spaces in Names

Spaces are used to separate the elements of the list that we are going to loop over. If one of those elements contains a space character, we need to surround it with quotes, and do the same thing to our loop variable. Suppose our data files are named:

```
red dragon.dat
purple unicorn.dat
```

To loop over these files, we would need to add double quotes like so:

```
$ for filename in "red dragon.dat" "purple unicorn.dat"
do
    head -n 100 "$filename" | tail -n 20
done
```

It is simpler just to avoid using spaces (or other special characters) in filenames.

The files above don't exist, so if we run the above code, the `head` command will be unable to find them, however the error message returned will show the name of the files it is expecting:

```
head: cannot open 'red dragon.dat' for reading: No such file or directory
head: cannot open 'purple unicorn.dat' for reading: No such file or directory
```

Try removing the quotes around `$filename` in the loop above to see the effect of the quote marks on spaces. Note that we get a result from the loop command for `unicorn.dat` when we run this code in the `creatures` directory:

```
head: cannot open 'red' for reading: No such file or directory
head: cannot open 'dragon.dat' for reading: No such file or directory
head: cannot open 'purple' for reading: No such file or directory
CGGTACCGAA
```



```
AAGGGTCGCG
CAAGTGTTCC
```

We would like to modify each of the files in `data-shell/creatures`, but also save a version of the original files, naming the copies `original-basilisk.dat` and `original-unicorn.dat`. We can't use:

```
$ cp *.dat original-*.dat
```

because that would expand to:

```
$ cp basilisk.dat unicorn.dat original-*.dat
```

This wouldn't back up our files, instead we get an error:

```
cp: target `original-*.dat' is not a directory
```

This problem arises when `cp` receives more than two inputs. When this happens, it expects the last input to be a directory where it can copy all the files it was passed. Since there is no directory named `original-*.dat` in the `creatures` directory we get an error.

Instead, we can use a loop:

```
$ for filename in *.dat
> do
>     cp $filename original-$filename
> done
```

This loop runs the `cp` command once for each filename. The first time, when `$filename` expands to `basilisk.dat`, the shell executes:

```
cp basilisk.dat original-basilisk.dat
```

The second time, the command is:

```
cp unicorn.dat original-unicorn.dat
```

Since the `cp` command does not normally produce any output, it's hard to check that the loop is doing the correct thing. However, we learned earlier how to print strings using `echo`, and we can modify the loop to use `echo` to print our commands without actually executing them. As such we can check what commands *would be* run in the unmodified loop.

The following diagram shows what happens when the modified loop is executed, and demonstrates how the judicious use of `echo` is a good debugging technique.

*TODO: Figure - For Loop in Action ../fig/shell\_script\_for\_loop\_flow\_chart.svg*

### Those Who Know History Can Choose to Repeat It

Another way to repeat previous work is to use the `history` command to get a list of the last few hundred commands that have been

executed, and then to use `!123` (where “123” is replaced by the command number) to repeat one of those commands. For example, if Beatrice types this:

```
$ history | tail -n 5
456  ls -l NENE0*.txt
457  rm stats-NENE01729B.txt.txt
458  bash goostats NENE01729B.txt stats-NENE01729B.txt
459  ls -l NENE0*.txt
460  history
```

then she can re-run `goostats` on `NENE01729B.txt` simply by typing `!458`.

### Other History Commands

There are a number of other shortcut commands for getting at the history.

- `Ctrl-R` enters a history search mode “reverse-i-search” and finds the most recent command in your history that matches the text you enter next. Press `Ctrl-R` one or more additional times to search for earlier matches.
- `!!` retrieves the immediately preceding command (you may or may not find this more convenient than using the up-arrow)
- `!$` retrieves the last word of the last command. That’s useful more often than you might expect: after `bash goostats NENE01729B.txt stats-NENE01729B.txt`, you can type `less !$` to look at the file `stats-NENE01729B.txt`, which is quicker than doing up-arrow and editing the command-line.

## 11.8 Shell Scripts

We are finally ready to see what makes the shell such a powerful programming environment. We are going to take the commands we repeat frequently and save them in files so that we can re-run all those operations again later by typing a single command. For historical reasons, a bunch of commands saved in a file is usually called a **shell script**, but make no mistake: these are actually small programs.

Let’s start by going back to `molecules/` and creating a new file, `middle.sh` which will become our shell script:

```
$ cd molecules
$ nano middle.sh
```

The command `nano middle.sh` opens the file `middle.sh` within the text editor “nano” (which runs within the shell). If the file does not exist, it will be created. We can use the text editor to directly edit the file – we’ll simply insert the following line:

```
head -n 15 octane.pdb | tail -n 5
```

This is a variation on the pipe we constructed earlier: it selects lines 11-15 of the file `octane.pdb`. Remember, we are *not* running it as a command just yet: we are putting the commands in a file.

Then we save the file (**Ctrl-O** in nano), and exit the text editor (**Ctrl-X** in nano). Check that the directory `molecules` now contains a file called `middle.sh`.

Once we have saved the file, we can ask the shell to execute the commands it contains. Our shell is called `bash`, so we run the following command:

```
$ bash middle.sh
```

```
ATOM      9  H           1      -4.502   0.681   0.785   1.00   0.00
ATOM     10  H           1      -5.254  -0.243  -0.537   1.00   0.00
ATOM     11  H           1      -4.357   1.252  -0.895   1.00   0.00
ATOM     12  H           1      -3.009  -0.741  -1.467   1.00   0.00
ATOM     13  H           1      -3.172  -1.337   0.206   1.00   0.00
```

Sure enough, our script’s output is exactly what we would get if we ran that pipeline directly.

### Text vs. Whatever

We usually call programs like Microsoft Word or LibreOffice Writer “text editors”, but we need to be a bit more careful when it comes to programming. By default, Microsoft Word uses `.docx` files to store not only text, but also formatting information about fonts, headings, and so on. This extra information isn’t stored as characters, and doesn’t mean anything to tools like `head`: they expect input files to contain nothing but the letters, digits, and punctuation on a standard computer keyboard. When editing programs, therefore, you must either use a plain text editor, or be careful to save files as plain text.

What if we want to select lines from an arbitrary file? We could edit `middle.sh` each time to change the filename, but that would probably take longer than just retyping the command. Instead, let’s edit `middle.sh` and make it more versatile:

```
$ nano middle.sh
```

Now, within “nano”, replace the text `octane.pdb` with the special variable called `$1`:

```
head -n 15 "$1" | tail -n 5
```

Inside a shell script, `$1` means “the first filename (or other argument) on the command line”. We can now run our script like this:

```
$ bash middle.sh octane.pdb
```

```
ATOM      9  H          1      -4.502   0.681   0.785   1.00   0.00
ATOM     10  H          1      -5.254  -0.243  -0.537   1.00   0.00
ATOM     11  H          1      -4.357   1.252  -0.895   1.00   0.00
ATOM     12  H          1      -3.009  -0.741  -1.467   1.00   0.00
ATOM     13  H          1      -3.172  -1.337   0.206   1.00   0.00
```

or on a different file like this:

```
$ bash middle.sh pentane.pdb
```

```
ATOM      9  H          1        1.324   0.350  -1.332   1.00   0.00
ATOM     10  H          1        1.271   1.378   0.122   1.00   0.00
ATOM     11  H          1       -0.074  -0.384   1.288   1.00   0.00
ATOM     12  H          1       -0.048  -1.362  -0.205   1.00   0.00
ATOM     13  H          1       -1.183   0.500  -1.412   1.00   0.00
```

### Double-Quotes Around Arguments

For the same reason that we put the loop variable inside double-quotes, in case the filename happens to contain any spaces, we surround `$1` with double-quotes.

We still need to edit `middle.sh` each time we want to adjust the range of lines, though. Let’s fix that by using the special variables `$2` and `$3` for the number of lines to be passed to `head` and `tail` respectively:

```
$ nano middle.sh
```

```
head -n "$2" "$1" | tail -n "$3"
```

We can now run:

```
$ bash middle.sh pentane.pdb 15 5
```

```
ATOM      9  H          1        1.324   0.350  -1.332   1.00   0.00
ATOM     10  H          1        1.271   1.378   0.122   1.00   0.00
ATOM     11  H          1       -0.074  -0.384   1.288   1.00   0.00
ATOM     12  H          1       -0.048  -1.362  -0.205   1.00   0.00
ATOM     13  H          1       -1.183   0.500  -1.412   1.00   0.00
```

By changing the arguments to our command we can change our script’s behaviour:

```
$ bash middle.sh pentane.pdb 20 5
```

```
ATOM     14  H          1       -1.259   1.420   0.112   1.00   0.00
ATOM     15  H          1       -2.608  -0.407   1.130   1.00   0.00
ATOM     16  H          1       -2.540  -1.303  -0.404   1.00   0.00
ATOM     17  H          1       -3.393   0.254  -0.321   1.00   0.00
```

```
TER          18          1
```

This works, but it may take the next person who reads `middle.sh` a moment to figure out what it does. We can improve our script by adding some **comments** at the top:

```
$ nano middle.sh

# Select lines from the middle of a file.
# Usage: bash middle.sh filename end_line num_lines
head -n "$2" "$1" | tail -n "$3"
```

A comment starts with a `#` character and runs to the end of the line. The computer ignores comments, but they're invaluable for helping people (including your future self) understand and use scripts. The only caveat is that each time you modify the script, you should check that the comment is still accurate: an explanation that sends the reader in the wrong direction is worse than none at all.

What if we want to process many files in a single pipeline? For example, if we want to sort our `.pdb` files by length, we would type:

```
$ wc -l *.pdb | sort -n
```

because `wc -l` lists the number of lines in the files (recall that `wc` stands for 'word count', adding the `-l` option means 'count lines' instead) and `sort -n` sorts things numerically. We could put this in a file, but then it would only ever sort a list of `.pdb` files in the current directory. If we want to be able to get a sorted list of other kinds of files, we need a way to get all those names into the script. We can't use `$1`, `$2`, and so on because we don't know how many files there are. Instead, we use the special variable `$@`, which means, "All of the command-line arguments to the shell script." We also should put `$@` inside double-quotes to handle the case of arguments containing spaces ("`$@`" is equivalent to "`$1`" "`$2`" ...) Here's an example:

```
$ nano sorted.sh

# Sort filenames by their length.
# Usage: bash sorted.sh one_or_more_filenames
wc -l "$@" | sort -n

$ bash sorted.sh *.pdb ../creatures/*.dat

9 methane.pdb
12 ethane.pdb
15 propane.pdb
20 cubane.pdb
21 pentane.pdb
30 octane.pdb
163 ../creatures/basilisk.dat
163 ../creatures/unicorn.dat
```

```
433 total
```

Suppose we have just run a series of commands that did something useful — for example, that created a graph we’d like to use in a paper. We’d like to be able to re-create the graph later if we need to, so we want to save the commands in a file. Instead of typing them in again (and potentially getting them wrong) we can do this:

```
$ history | tail -n 5 > redo-figure-3.sh
```

The file `redo-figure-3.sh` now contains:

```
297 bash goostats NENE01729B.txt stats-NENE01729B.txt
298 bash goodiff stats-NENE01729B.txt /data/validated/01729.txt > 01729-differences.txt
299 cut -d ',' -f 2-3 01729-differences.txt > 01729-time-series.txt
300 ygraph --format scatter --color bw --borders none 01729-time-series.txt figure-3.p
301 history | tail -n 5 > redo-figure-3.sh
```

After a moment’s work in an editor to remove the serial numbers on the commands, and to remove the final line where we called the `history` command, we have a completely accurate record of how we created that figure.

In practice, most people develop shell scripts by running commands at the shell prompt a few times to make sure they’re doing the right thing, then saving them in a file for re-use. This style of work allows people to recycle what they discover about their data and their workflow with one call to `history` and a bit of editing to clean up the output and save it as a shell script.

## 11.9 Finding Things

In the same way that many of us now use “Google” as a verb meaning “to find”, Unix programmers often use the word “grep”. “grep” is a contraction of “global/regular expression/print”, a common sequence of operations in early Unix text editors. It is also the name of a very useful command-line program.

`grep` finds and prints lines in files that match a pattern. For our examples, we will use a file that contains three haikus taken from a 1998 competition in *Salon* magazine. For this set of examples, we’re going to be working in the writing subdirectory:

```
$ cd
$ cd Desktop/data-shell/writing
$ cat haiku.txt
```

```
The Tao that is seen
Is not the true Tao, until
You bring fresh toner.
```

```
With searching comes loss
and the presence of absence:
"My Thesis" not found.
```

```
Yesterday it worked
Today it is not working
Software is like that.
```

### Forever, or Five Years

We haven't linked to the original haikus because they don't appear to be on *Salon's* site any longer. As Jeff Rothenberg said, "Digital information lasts forever — or five years, whichever comes first." Luckily, popular content often has backups.

Let's find lines that contain the word "not":

```
$ grep not haiku.txt

Is not the true Tao, until
"My Thesis" not found
Today it is not working
```

Here, **not** is the pattern we're searching for. The **grep** command searches through the file, looking for matches to the pattern specified. To use it type **grep**, then the pattern we're searching for and finally the name of the file (or files) we're searching in.

The output is the three lines in the file that contain the letters "not".

By default, **grep** searches for a pattern in a case-sensitive way. In addition, the search pattern we have selected does not have to form a complete word, as we will see in the next example.

Let's search for the pattern: "The".

```
$ grep The haiku.txt

The Tao that is seen
"My Thesis" not found.
```

This time, two lines that include the letters "The" are outputted, one of which contained our search pattern within a larger word, "Thesis".

To restrict matches to lines containing the word "The" on its own, we can give **grep** with the **-w** option. This will limit matches to word boundaries.

Later in this lesson, we will also see how we can change the search behavior of **grep** with respect to its case sensitivity.

```
$ grep -w The haiku.txt

The Tao that is seen
```

Note that a “word boundary” includes the start and end of a line, so not just letters surrounded by spaces. Sometimes we don’t want to search for a single word, but a phrase. This is also easy to do with `grep` by putting the phrase in quotes.

```
$ grep -w "is not" haiku.txt
```

```
Today it is not working
```

We’ve now seen that you don’t have to have quotes around single words, but it is useful to use quotes when searching for multiple words. It also helps to make it easier to distinguish between the search term or phrase and the file being searched. We will use quotes in the remaining examples.

Another useful option is `-n`, which numbers the lines that match:

```
$ grep -n "it" haiku.txt
```

```
5:With searching comes loss
9:Yesterday it worked
10:Today it is not working
```

Here, we can see that lines 5, 9, and 10 contain the letters “it”.

We can combine options (i.e. flags) as we do with other Unix commands. For example, let’s find the lines that contain the word “the”. We can combine the option `-w` to find the lines that contain the word “the” and `-n` to number the lines that match:

```
$ grep -n -w "the" haiku.txt
```

```
2:Is not the true Tao, until
6:and the presence of absence:
```

Now we want to use the option `-i` to make our search case-insensitive:

```
$ grep -n -w -i "the" haiku.txt
```

```
1:The Tao that is seen
2:Is not the true Tao, until
6:and the presence of absence:
```

Now, we want to use the option `-v` to invert our search, i.e., we want to output the lines that do not contain the word “the”.

```
$ grep -n -w -v "the" haiku.txt
```

```
1:The Tao that is seen
3:You bring fresh toner.
4:
5:With searching comes loss
7:"My Thesis" not found.
8:
9:Yesterday it worked
```



```
10:Today it is not working
11:Software is like that.
```

`grep` has lots of other options. To find out what they are, we can type:

```
$ grep --help
```

```
Usage: grep [OPTION]... PATTERN [FILE]...
Search for PATTERN in each FILE or standard input.
PATTERN is, by default, a basic regular expression (BRE).
Example: grep -i 'hello world' menu.h main.c
```

Regex selection and interpretation:

|                                    |   |
|------------------------------------|---|
| <code>-E, --extended-regexp</code> | PATTERN is an extended regular expression (ERE)     |
| <code>-F, --fixed-strings</code>   | PATTERN is a set of newline-separated fixed strings |
| <code>-G, --basic-regexp</code>    | PATTERN is a basic regular expression (BRE)         |
| <code>-P, --perl-regexp</code>     | PATTERN is a Perl regular expression                |
| <code>-e, --regexp=PATTERN</code>  | use PATTERN for matching                            |
| <code>-f, --file=FILE</code>       | obtain PATTERN from FILE                            |
| <code>-i, --ignore-case</code>     | ignore case distinctions                            |
| <code>-w, --word-regexp</code>     | force PATTERN to match only whole words             |
| <code>-x, --line-regexp</code>     | force PATTERN to match only whole lines             |
| <code>-z, --null-data</code>       | a data line ends in 0 byte, not newline             |

Miscellaneous:

```
...      ...      ...
```

### Wildcards

`grep`'s real power doesn't come from its options, though; it comes from the fact that patterns can include wildcards. (The technical name for these is **regular expressions**, which is what the "re" in "grep" stands for.) Regular expressions are both complex and powerful; if you want to do complex searches, please look at the lesson on our website. As a taster, we can find lines that have an 'o' in the second position like this:

```
$ grep -E '^..o' haiku.txt
```

```
You bring fresh toner.
Today it is not working
Software is like that.
```

We use the `-E` option and put the pattern in quotes to prevent the shell from trying to interpret it. (If the pattern contained a `*`, for example, the shell would try to expand it before running `grep`.) The `^` in the pattern anchors the match to the start of the line. The `.` matches a single character (just like `?` in the shell), while the `o` matches an actual 'o'.

While **grep** finds lines in files, the **find** command finds files themselves. Again, it has a lot of options; to show how the simplest ones work, we'll use the directory tree shown below.

*TODO: Figure - File Tree for Find Example ../fig/find-file-tree.svg*

Beatrice's **writing** directory contains one file called **haiku.txt** and three sub-directories: **thesis** (which contains a sadly empty file, **empty-draft.md**); **data** (which contains three files **LittleWomen.txt**, **one.txt** and **two.txt**); and a **tools** directory that contains the programs **format** and **stats**, and a subdirectory called **old**, with a file **oldtool**.

For our first command, let's run **find .** (remember to run this command from the **data-shell/writing** folder).

```
$ find .
.
./data
./data/one.txt
./data/LittleWomen.txt
./data/two.txt
./tools
./tools/format
./tools/old
./tools/old/oldtool
./tools/stats
./haiku.txt
./thesis
./thesis/empty-draft.md
```

As always, the **.** on its own means the current working directory, which is where we want our search to start. **find**'s output is the names of every file **and** directory under the current working directory. This can seem useless at first but **find** has many options to filter the output and in this lesson we will discover some of them.

The first option in our list is **-type d** that means "things that are directories". Sure enough, **find**'s output is the names of the five directories in our little tree (including **.**):

```
$ find . -type d
./
./data
./thesis
./tools
./tools/old
```

Notice that the objects **find** finds are not listed in any particular order. If we change **-type d** to **-type f**, we get a listing of all the files instead:

```
$ find . -type f
./haiku.txt
./tools/stats
./tools/old/oldtool
./tools/format
./thesis/empty-draft.md
./data/one.txt
./data/LittleWomen.txt
./data/two.txt
```

Now let's try matching by name:

```
$ find . -name *.txt
./haiku.txt
```

We expected it to find all the text files, but it only prints out `./haiku.txt`. The problem is that the shell expands wildcard characters like `*` *before* commands run. Since `*.txt` in the current directory expands to `haiku.txt`, the command we actually ran was:

```
$ find . -name haiku.txt
```

`find` did what we asked; we just asked for the wrong thing.

To get what we want, let's do what we did with `grep`: put `*.txt` in single quotes to prevent the shell from expanding the `*` wildcard. This way, `find` actually gets the pattern `*.txt`, not the expanded filename `haiku.txt`:

```
$ find . -name '*.txt'
./data/one.txt
./data/LittleWomen.txt
./data/two.txt
./haiku.txt
```

### Listing vs. Finding

`ls` and `find` can be made to do similar things given the right options, but under normal circumstances, `ls` lists everything it can, while `find` searches for things with certain properties and shows them.

As we said earlier, the command line's power lies in combining tools. We've seen how to do that with pipes; let's look at another technique. As we just saw, `find . -name '*.txt'` gives us a list of all text files in or below the current directory. How can we combine that with `wc -l` to count the lines in all those files?

The simplest way is to put the `find` command inside `$()`:

```
$ wc -l $(find . -name '*.txt')
```

```
11 ./haiku.txt
300 ./data/two.txt
21022 ./data/LittleWomen.txt
70 ./data/one.txt
21403 total
```

When the shell executes this command, the first thing it does is run whatever is inside the `$()`. It then replaces the `$()` expression with that command's output. Since the output of `find` is the four filenames `./data/one.txt`, `./data/LittleWomen.txt`, `./data/two.txt`, and `./haiku.txt`, the shell constructs the command:

```
$ wc -l ./data/one.txt ./data/LittleWomen.txt ./data/two.txt ./haiku.txt
```

which is what we wanted. This expansion is exactly what the shell does when it expands wildcards like `*` and `?`, but lets us use any command we want as our own “wildcard”.

It's very common to use `find` and `grep` together. The first finds files that match a pattern; the second looks for lines inside those files that match another pattern. Here, for example, we can find PDB files that contain iron atoms by looking for the string “FE” in all the `.pdb` files above the current directory:

```
$ grep "FE" $(find .. -name '*.pdb')

../data/pdb/heme.pdb:ATOM      25  FE                      1      -0.924   0.535  -0.518
```

### Binary Files

We have focused exclusively on finding patterns in text files. What if your data is stored as images, in databases, or in some other format?

A handful of tools extend `grep` to handle a few non text formats. But a more generalizable approach is to convert the data to text, or extract the text-like elements from the data. On the one hand, it makes simple things easy to do. On the other hand, complex things are usually impossible. For example, it's easy enough to write a program that will extract X and Y dimensions from image files for `grep` to play with, but how would you write something to find values in a spreadsheet whose cells contained formulas?

A last option is to recognize that the shell and text processing have their limits, and to use another programming language. When the time comes to do this, don't be too hard on the shell: many modern programming languages have borrowed a lot of ideas from it, and imitation is also the sincerest form of praise.

## 11.10 Summary

What have we learned in these lessons on the bash shell? The Navigating Files and Directories and Working With Files and Directories content means that we can now get around in a command line environment. This is especially useful if we are doing our data analysis by logging onto a high performance computing facility, for instance, as the command line interface is often the only way to interact with such a facility.

The Loops, Shell Scripts and Finding Things content means that we can now save time by getting the computer to do repetitive tasks for us.

Finally, the Pipes and Filters content means we've been introduced to a fundamentally important concept: rather than writing monolithic functions/programs that do many things at once, it's better to write small functions/programs that do perform a well-defined, discrete task, and then combine them in myriad different ways to conduct your data analysis.

## 11.11 Exercises

### 11.11.1 Navigating Files and Directories

#### Exploring More `ls` Flags

You can also use two options at the same time. What does the command `ls` do when used with the `-l` option? What about if you use both the `-l` and the `-h` option?

Some of its output is about properties that we do not cover in this lesson (such as file permissions and ownership), but the rest should be useful nevertheless.

#### *Solution*

The `-l` option makes `ls` use a long listing format, showing not only the file/directory names but also additional information such as the file size and the time of its last modification. If you use both the `-h` option and the `-l` option, this makes the file size “human readable”, i.e. displaying something like 5.3K instead of 5369.

#### Listing Recursively and By Time

The command `ls -R` lists the contents of directories recursively, i.e., lists their sub-directories, sub-sub-directories, and so on at each level. The command `ls -t` lists things by time of last change, with most recently changed files or directories first.

In what order does `ls -R -t` display things? Hint: `ls -l` uses a long listing format to view timestamps.

*Solution*

The files/directories in each directory are sorted by time of last change.

**Absolute vs Relative Paths**

Starting from `/Users/amanda/data`, which of the following commands could Amanda use to navigate to her home directory, which is `/Users/amanda`?

1. `cd .`
2. `cd /`
3. `cd /home/amanda`
4. `cd ../../`
5. `cd ~`
6. `cd home`
7. `cd ~/data/..`
8. `cd`
9. `cd ..`

*Solution*

1. No: `.` stands for the current directory.
2. No: `/` stands for the root directory.
3. No: Amanda's home directory is `/Users/amanda`.
4. No: this goes up two levels, i.e. ends in `/Users`.
5. Yes: `~` stands for the user's home directory, in this case `/Users/amanda`.
6. No: this would navigate into a directory `home` in the current directory if it exists.
7. Yes: unnecessarily complicated, but correct.
8. Yes: shortcut to go back to the user's home directory.
9. Yes: goes up one level.

**Relative Path Resolution**

Using the filesystem diagram below, if `pwd` displays `/Users/thing`, what will `ls -F ../backup` display?

1. `../backup: No such file or directory`
2. `2012-12-01 2013-01-08 2013-01-27`
3. `2012-12-01/ 2013-01-08/ 2013-01-27/`
4. `original/ pnas_final/ pnas_sub/`

TODO: Figure - File System for Challenge Questions `../fig/filesystem-challenge.svg`

*Solution*

1. No: there *is* a directory `backup` in `/Users`.
2. No: this is the content of `Users/thing/backup`, but with `..` we asked for one level further up.

3. No: see previous explanation.
4. Yes: `../backup/` refers to `/Users/backup/`.

### ls Reading Comprehension

Using the filesystem diagram below, if `pwd` displays `/Users/backup`, and `-r` tells `ls` to display things in reverse order, what command(s) will result in the following output:

```
pnas_sub/ pnas_final/ original/
```

TODO: Figure - File System for Challenge Questions `../fig/filesystem-challenge.svg`

1. `ls pwd`
2. `ls -r -F`
3. `ls -r -F /Users/backup`

#### *Solution*

1. No: `pwd` is not the name of a directory.
2. Yes: `ls` without directory argument lists files and directories in the current directory.
3. Yes: uses the absolute path explicitly.

## 11.11.2 Working With Files and Directories

### Creating Files a Different Way

We have seen how to create text files using the `nano` editor. Now, try the following command:

```
$ touch my_file.txt
```

1. What did the `touch` command do? When you look at your current directory using the GUI file explorer, does the file show up?
2. Use `ls -l` to inspect the files. How large is `my_file.txt`?
3. When might you want to create a file this way?

#### *Solution*

1. The `touch` command generates a new file called `my_file.txt` in your current directory. You can observe this newly generated file by typing `ls` at the command line prompt. `my_file.txt` can also be viewed in your GUI file explorer.
2. When you inspect the file with `ls -l`, note that the size of `my_file.txt` is 0 bytes. In other words, it contains no data. If you open `my_file.txt` using your text editor it is blank.

3. Some programs do not generate output files themselves, but instead require that empty files have already been generated. When the program is run, it searches for an existing file to populate with its output. The `touch` command allows you to efficiently generate a blank text file to be used by such programs.

### Moving to the Current Folder

After running the following commands, Jamie realizes that she put the files `sucrose.dat` and `maltose.dat` into the wrong folder:

```
$ ls -F
  analyzed/ raw/
$ ls -F analyzed
  fructose.dat glucose.dat maltose.dat sucrose.dat
$ cd raw/
```

Fill in the blanks to move these files to the current folder (i.e., the one she is currently in):

```
$ mv ___/sucrose.dat ___/maltose.dat ___
```

*Solution*

```
$ mv ../analyzed/sucrose.dat ../analyzed/maltose.dat .
```

Recall that `..` refers to the parent directory (i.e. one above the current directory) and that `.` refers to the current directory.

### Renaming Files

Suppose that you created a plain-text file in your current directory to contain a list of the statistical tests you will need to do to analyze your data, and named it: `statstics.txt`

After creating and saving this file you realize you misspelled the filename! You want to correct the mistake, which of the following commands could you use to do so?

1. `cp statstics.txt statistics.txt`
2. `mv statstics.txt statistics.txt`
3. `mv statstics.txt .`
4. `cp statstics.txt .`

*Solution*

1. No. While this would create a file with the correct name, the incorrectly named file still exists in the directory and would need to be deleted.
2. Yes, this would work to rename the file.
3. No, the period(`.`) indicates where to move the file, but does not provide a new file name; identical file names cannot be created.



4. No, the period(.) indicates where to copy the file, but does not provide a new file name; identical file names cannot be created.

### Moving and Copying

What is the output of the closing `ls` command in the sequence shown below?

```
$ pwd
/Users/jamie/data
$ ls
proteins.dat
$ mkdir recombine
$ mv proteins.dat recombine/
$ cp recombine/proteins.dat ../proteins-saved.dat
$ ls
1. proteins-saved.dat recombine
2. recombine
3. proteins.dat recombine
4. proteins-saved.dat
```

#### *Solution*

We start in the `/Users/jamie/data` directory, and create a new folder called `recombine`. The second line moves (`mv`) the file `proteins.dat` to the new folder (`recombine`). The third line makes a copy of the file we just moved. The tricky part here is where the file was copied to. Recall that `..` means “go up a level”, so the copied file is now in `/Users/jamie`. Notice that `..` is interpreted with respect to the current working directory, **not** with respect to the location of the file being copied. So, the only thing that will show using `ls` (in `/Users/jamie/data`) is the `recombine` folder.

1. No, see explanation above. `proteins-saved.dat` is located at `/Users/jamie`
2. Yes
3. No, see explanation above. `proteins.dat` is located at `/Users/jamie/data/recombine`
4. No, see explanation above. `proteins-saved.dat` is located at `/Users/jamie`

### Using `rm` Safely

What happens when we execute `rm -i thesis_backup/quotations.txt`? Why would we want this protection when using `rm`?

#### *Solution*

```
$ rm: remove regular file 'thesis_backup/quotations.txt'? y
```

The `-i` option will prompt before (every) removal (use Y to confirm deletion or N to keep the file). The Unix shell doesn't have a trash bin, so all the files removed will disappear forever. By using the `-i` option, we have the chance to check that we are deleting only the files that we want to remove.

### Copy with Multiple Filenames

For this exercise, you can test the commands in the `data-shell/data` directory.

In the example below, what does `cp` do when given several filenames and a directory name?

```
$ mkdir backup
$ cp amino-acids.txt animals.txt backup/
```

In the example below, what does `cp` do when given three or more file names?

```
$ ls -F
amino-acids.txt  animals.txt  backup/  elements/  morse.txt  pdb/  planets.txt  salmon
$ cp amino-acids.txt animals.txt morse.txt
```

#### *Solution*

If given more than one file name followed by a directory name (i.e. the destination directory must be the last argument), `cp` copies the files to the named directory.

If given three file names, `cp` throws an error such as the one below, because it is expecting a directory name as the last argument.

```
cp: target 'morse.txt' is not a directory
```

### List filenames matching a pattern

When run in the `molecules` directory, which `ls` command(s) will produce this output?

```
ethane.pdb  methane.pdb
```

1. `ls *t*ane.pdb`
2. `ls *t?ne.*`
3. `ls *t??ne.pdb`
4. `ls ethane.*`

#### *Solution*

The solution is 3.

1. shows all files whose names contain zero or more characters (\*) followed by the letter `t`, then zero or more characters (\*) followed by `ane.pdb`. This gives `ethane.pdb` `methane.pdb` `octane.pdb` `pentane.pdb`.

2. shows all files whose names start with zero or more characters (\*) followed by the letter `t`, then a single character (?), then `ne.` followed by zero or more characters (\*). This will give us `octane.pdb` and `pentane.pdb` but doesn't match anything which ends in `thane.pdb`.
3. fixes the problems of option 2 by matching two characters (??) between `t` and `ne`. This is the solution.
4. only shows files starting with `ethane.`

### More on Wildcards

Sam has a directory containing calibration data, datasets, and descriptions of the datasets:

```
.
  2015-10-23-calibration.txt
  2015-10-23-dataset1.txt
  2015-10-23-dataset2.txt
  2015-10-23-dataset_overview.txt
  2015-10-26-calibration.txt
  2015-10-26-dataset1.txt
  2015-10-26-dataset2.txt
  2015-10-26-dataset_overview.txt
  2015-11-23-calibration.txt
  2015-11-23-dataset1.txt
  2015-11-23-dataset2.txt
  2015-11-23-dataset_overview.txt
  backup
    calibration
    datasets
  send_to_bob
    all_datasets_created_on_a_23rd
    all_november_files
```

Before heading off to another field trip, she wants to back up her data and send some datasets to her colleague Bob. Sam uses the following commands to get the job done:

```
$ cp *dataset* backup/datasets
$ cp ____calibration____ backup/calibration
$ cp 2015-____-____ send_to_bob/all_november_files/
$ cp ____ send_to_bob/all_datasets_created_on_a_23rd/
```

Help Sam by filling in the blanks.

The resulting directory structure should look like this

```
.
  2015-10-23-calibration.txt
```

```

2015-10-23-dataset1.txt
2015-10-23-dataset2.txt
2015-10-23-dataset_overview.txt
2015-10-26-calibration.txt
2015-10-26-dataset1.txt
2015-10-26-dataset2.txt
2015-10-26-dataset_overview.txt
2015-11-23-calibration.txt
2015-11-23-dataset1.txt
2015-11-23-dataset2.txt
2015-11-23-dataset_overview.txt
backup
  calibration
    2015-10-23-calibration.txt
    2015-10-26-calibration.txt
    2015-11-23-calibration.txt
  datasets
    2015-10-23-dataset1.txt
    2015-10-23-dataset2.txt
    2015-10-23-dataset_overview.txt
    2015-10-26-dataset1.txt
    2015-10-26-dataset2.txt
    2015-10-26-dataset_overview.txt
    2015-11-23-dataset1.txt
    2015-11-23-dataset2.txt
    2015-11-23-dataset_overview.txt
send_to_bob
  all_datasets_created_on_a_23rd
    2015-10-23-dataset1.txt
    2015-10-23-dataset2.txt
    2015-10-23-dataset_overview.txt
    2015-11-23-dataset1.txt
    2015-11-23-dataset2.txt
    2015-11-23-dataset_overview.txt
  all_november_files
    2015-11-23-calibration.txt
    2015-11-23-dataset1.txt
    2015-11-23-dataset2.txt
    2015-11-23-dataset_overview.txt

```

*Solution*

```

$ cp *calibration.txt backup/calibration
$ cp 2015-11-* send_to_bob/all_november_files/
$ cp *-23-dataset* send_to_bob/all_datasets_created_on_a_23rd/

```

**Organizing Directories and Files**

Jamie is working on a project and she sees that her files aren't very well organized:

```
$ ls -F
analyzed/  fructose.dat  raw/  sucrose.dat
```

The `fructose.dat` and `sucrose.dat` files contain output from her data analysis. What command(s) covered in this lesson does she need to run so that the commands below will produce the output shown?

```
$ ls -F
analyzed/  raw/
$ ls analyzed
fructose.dat  sucrose.dat
```

*Solution*

```
mv *.dat analyzed
```

Jamie needs to move her files `fructose.dat` and `sucrose.dat` to the `analyzed` directory. The shell will expand `*.dat` to match all `.dat` files in the current directory. The `mv` command then moves the list of `.dat` files to the “analyzed” directory.

### Reproduce a folder structure

You're starting a new experiment, and would like to duplicate the directory structure from your previous experiment so you can add new data.

Assume that the previous experiment is in a folder called '2016-05-18', which contains a `data` folder that in turn contains folders named `raw` and `processed` that contain data files. The goal is to copy the folder structure of the `2016-05-18-data` folder into a folder called `2016-05-20` so that your final directory structure looks like this:

```
2016-05-20/
  data
    processed
    raw
```

Which of the following set of commands would achieve this objective?

What would the other commands do?

```
$ mkdir 2016-05-20
$ mkdir 2016-05-20/data
$ mkdir 2016-05-20/data/processed
$ mkdir 2016-05-20/data/raw

$ mkdir 2016-05-20
$ cd 2016-05-20
```

```
$ mkdir data
$ cd data
$ mkdir raw processed

$ mkdir 2016-05-20/data/raw
$ mkdir 2016-05-20/data/processed

$ mkdir 2016-05-20
$ cd 2016-05-20
$ mkdir data
$ mkdir raw processed
```

#### *Solution*

The first two sets of commands achieve this objective. The first set uses relative paths to create the top level directory before the subdirectories.

The third set of commands will give an error because `mkdir` won't create a subdirectory of a non-existent directory: the intermediate level folders must be created first.

The final set of commands generates the 'raw' and 'processed' directories at the same level as the 'data' directory.

### 11.11.3 Pipes and Filters

#### What Does `sort -n` Do?

If we run `sort` on a file containing the following lines:

```
10
2
19
22
6
```

the output is:

```
10
19
2
22
6
```

If we run `sort -n` on the same input, we get this instead:

```
2
6
10
19
```

22

Explain why `-n` has this effect.

*Solution*

The `-n` option specifies a numerical rather than an alphanumerical sort.

### What Does `>>` Mean?

We have seen the use of `>`, but there is a similar operator `>>` which works slightly differently. We'll learn about the differences between these two operators by printing some strings. We can use the `echo` command to print strings e.g.

```
$ echo The echo command prints text
```

The echo command prints text

Now test the commands below to reveal the difference between the two operators:

```
$ echo hello > testfile01.txt
```

and:

```
$ echo hello >> testfile02.txt
```

Hint: Try executing each command twice in a row and then examining the output files.

*Solution*

In the first example with `>`, the string "hello" is written to `testfile01.txt`, but the file gets overwritten each time we run the command.

We see from the second example that the `>>` operator also writes "hello" to a file (in this case `testfile02.txt`), but appends the string to the file if it already exists (i.e. when we run it for the second time).

### Appending Data

We have already met the `head` command, which prints lines from the start of a file. `tail` is similar, but prints lines from the end of a file instead.

Consider the file `data-shell/data/animals.txt`. After these commands, select the answer that corresponds to the file `animals-subset.txt`:

```
$ head -n 3 animals.txt > animals-subset.txt
```

```
$ tail -n 2 animals.txt >> animals-subset.txt
```

1. The first three lines of `animals.txt`
2. The last two lines of `animals.txt`
3. The first three lines and the last two lines of `animals.txt`
4. The second and third lines of `animals.txt`

*Solution* Option 3 is correct. For option 1 to be correct we would only run the `head` command. For option 2 to be correct we would only run the `tail` command. For option 4 to be correct we would have to pipe the output of `head` into `tail -n 2` by doing `head -n 3 animals.txt | tail -n 2 > animals-subset.txt`

### Piping Commands Together

In our current directory, we want to find the 3 files which have the least number of lines. Which command listed below would work?

1. `wc -l * > sort -n > head -n 3`
2. `wc -l * | sort -n | head -n 1-3`
3. `wc -l * | head -n 3 | sort -n`
4. `wc -l * | sort -n | head -n 3`

*Solution*

Option 4 is the solution. The pipe character `|` is used to feed the standard output from one process to the standard input of another. `>` is used to redirect standard output to a file. Try it in the `data-shell/molecules` directory!

### Why Does `uniq` Only Remove Adjacent Duplicates?

The command `uniq` removes adjacent duplicated lines from its input. For example, the file `data-shell/data/salmon.txt` contains:

```
coho
coho
steelhead
coho
steelhead
steelhead
```

Running the command `uniq salmon.txt` from the `data-shell/data` directory produces:

```
coho
steelhead
coho
steelhead
```

Why do you think `uniq` only removes *adjacent* duplicated lines? (Hint: think about very large data sets.) What other command could you combine with it in a pipe to remove all duplicated lines?

*Solution*

```
$ sort salmon.txt | uniq
```

### Pipe Reading Comprehension



A file called `animals.txt` (in the `data-shell/data` folder) contains the following data:

```
2012-11-05,deer
2012-11-05,rabbit
2012-11-05,raccoon
2012-11-06,rabbit
2012-11-06,deer
2012-11-06,fox
2012-11-07,rabbit
2012-11-07,bear
```

What text passes through each of the pipes and the final redirect in the pipeline below?

```
$ cat animals.txt | head -n 5 | tail -n 3 | sort -r > final.txt
```

Hint: build the pipeline up one command at a time to test your understanding

#### *Solution*

The `head` command extracts the first 5 lines from `animals.txt`. Then, the last 3 lines are extracted from the previous 5 by using the `tail` command. With the `sort -r` command those 3 lines are sorted in reverse order and finally, the output is redirected to a file `final.txt`. The content of this file can be checked by executing `cat final.txt`. The file should contain the following lines:

```
2012-11-06,rabbit
2012-11-06,deer
2012-11-05,raccoon
```

#### **Pipe Construction**

For the file `animals.txt` from the previous exercise, consider the following command:

```
$ cut -d , -f 2 animals.txt
```

The `cut` command is used to remove or “cut out” certain sections of each line in the file. The `-d` option is used to define the delimiter. A **delimiter** is a character that is used to separate each line of text into columns. The default delimiter is Tab, meaning that the `cut` command will automatically assume that values in different columns will be separated by a tab. The `-f` option is used to specify the field (column) to cut out. The command above uses the `-d` option to split each line by comma, and the `-f` option to print the second field in each line, to give the following output:

```
deer
rabbit
raccoon
rabbit
```

```

deer
fox
rabbit
bear

```

What other command(s) could be added to this in a pipeline to find out what animals the file contains (without any duplicates in their names)?

*Solution*

```
$ cut -d , -f 2 animals.txt | sort | uniq
```

### Which Pipe?

The file `animals.txt` contains 8 lines of data formatted as follows:

```

2012-11-05,deer
2012-11-05,rabbit
2012-11-05,raccoon
2012-11-06,rabbit
...

```

The `uniq` command has a `-c` option which gives a count of the number of times a line occurs in its input. Assuming your current directory is `data-shell/data/`, what command would you use to produce a table that shows the total count of each type of animal in the file?

1. `sort animals.txt | uniq -c`
2. `sort -t, -k2,2 animals.txt | uniq -c`
3. `cut -d, -f 2 animals.txt | uniq -c`
4. `cut -d, -f 2 animals.txt | sort | uniq -c`
5. `cut -d, -f 2 animals.txt | sort | uniq -c | wc -l`

*Solution*

Option 4. is the correct answer. If you have difficulty understanding why, try running the commands, or sub-sections of the pipelines (make sure you are in the `data-shell/data` directory).

### Wildcard Expressions

Wildcard expressions can be very complex, but you can sometimes write them in ways that only use simple syntax, at the expense of being a bit more verbose. Consider the directory `data-shell/north-pacific-gyre/2012-07-03` : the wildcard expression `*[AB].txt` matches all files ending in `A.txt` or `B.txt`. Imagine you forgot about this.

1. Can you match the same set of files with basic wildcard expressions that do not use the `[]` syntax? *Hint*: You may need more than one expression.
2. The expression that you found and the expression from the lesson match the same set of files in this example. What is the small difference between the outputs?

3. Under what circumstances would your new expression produce an error message where the original one would not?

*Solution*

1. A solution using two wildcard expressions: `shell $ ls *A.txt $ ls *B.txt`
2. The output from the new commands is separated because there are two commands.
3. When there are no files ending in `A.txt`, or there are no files ending in `B.txt`.

### Removing Unneeded Files

Suppose you want to delete your processed data files, and only keep your raw files and processing script to save storage. The raw files end in `.dat` and the processed files end in `.txt`. Which of the following would remove all the processed data files, and *only* the processed data files?

1. `rm ?.txt`
2. `rm *.txt`
3. `rm * .txt`
4. `rm *.*`

*Solution*

1. This would remove `.txt` files with one-character names
2. This is correct answer
3. The shell would expand `*` to match everything in the current directory, so the command would try to remove all matched files and an additional file called `.txt`
4. The shell would expand `*.*` to match all files with any extension, so this command would delete all files

#### 11.11.4 Loops

##### Doing a Dry Run

A loop is a way to do many things at once — or to make many mistakes at once if it does the wrong thing. One way to check what a loop *would* do is to `echo` the commands it would run instead of actually running them.

Suppose we want to preview the commands the following loop will execute without actually running those commands:

```
$ for file in *.pdb
> do
>   analyze $file > analyzed-$file
> done
```

What is the difference between the two loops below, and which one would we want to run?

```
# Version 1
$ for file in *.pdb
> do
>   echo analyze $file > analyzed-$file
> done

# Version 2
$ for file in *.pdb
> do
>   echo "analyze $file > analyzed-$file"
> done
```

#### *Solution*

The second version is the one we want to run. This prints to screen everything enclosed in the quote marks, expanding the loop variable name because we have prefixed it with a dollar sign.

The first version redirects the output from the command `echo analyze $file` to a file, `analyzed-$file`. A series of files is generated: `analyzed-cubane.pdb`, `analyzed-ethane.pdb` etc.

Try both versions for yourself to see the output! Be sure to open the `analyzed-*.pdb` files to view their contents.

### **Nested Loops**

Suppose we want to set up a directory structure to organize some experiments measuring reaction rate constants with different compounds *and* different temperatures. What would be the result of the following code:

```
$ for species in cubane ethane methane
> do
>   for temperature in 25 30 37 40
>   do
>     mkdir $species-$temperature
>   done
> done
```

#### *Solution*

We have a nested loop, i.e. contained within another loop, so for each species in the outer loop, the inner loop (the nested loop) iterates over the list of temperatures, and creates a new directory for each combination.

Try running the code for yourself to see which directories are created!

### **Variables in Loops**

This exercise refers to the `data-shell/molecules` directory. `ls` gives the following output:

```
cubane.pdb  ethane.pdb  methane.pdb  octane.pdb  pentane.pdb  propane.pdb
```

What is the output of the following code?

```
$ for datafile in *.pdb
> do
>   ls *.pdb
> done
```

Now, what is the output of the following code?

```
$ for datafile in *.pdb
> do
>   ls $datafile
> done
```

Why do these two loops give different outputs?

#### *Solution*

The first code block gives the same output on each iteration through the loop. Bash expands the wildcard `*.pdb` within the loop body (as well as before the loop starts) to match all files ending in `.pdb` and then lists them using `ls`. The expanded loop would look like this:

```
$ for datafile in cubane.pdb  ethane.pdb  methane.pdb  octane.pdb  pentane.pdb  propane.pdb
> do
>   ls cubane.pdb  ethane.pdb  methane.pdb  octane.pdb  pentane.pdb  propane.pdb
> done

cubane.pdb  ethane.pdb  methane.pdb  octane.pdb  pentane.pdb  propane.pdb
cubane.pdb  ethane.pdb  methane.pdb  octane.pdb  pentane.pdb  propane.pdb
cubane.pdb  ethane.pdb  methane.pdb  octane.pdb  pentane.pdb  propane.pdb
cubane.pdb  ethane.pdb  methane.pdb  octane.pdb  pentane.pdb  propane.pdb
cubane.pdb  ethane.pdb  methane.pdb  octane.pdb  pentane.pdb  propane.pdb
cubane.pdb  ethane.pdb  methane.pdb  octane.pdb  pentane.pdb  propane.pdb
```

The second code block lists a different file on each loop iteration. The value of the `datafile` variable is evaluated using `$datafile`, and then listed using `ls`.

```
cubane.pdb
ethane.pdb
methane.pdb
octane.pdb
pentane.pdb
propane.pdb
```

### Limiting Sets of Files

What would be the output of running the following loop in the `data-shell/molecules` directory?

```
$ for filename in c*
> do
>   ls $filename
> done
```

1. No files are listed.
2. All files are listed.
3. Only `cubane.pdb`, `octane.pdb` and `pentane.pdb` are listed.
4. Only `cubane.pdb` is listed.

#### *Solution*

4 is the correct answer. `*` matches zero or more characters, so any file name starting with the letter `c`, followed by zero or more other characters will be matched.

How would the output differ from using this command instead?

```
$ for filename in *c*
> do
>   ls $filename
> done
```

1. The same files would be listed.
2. All the files are listed this time.
3. No files are listed this time.
4. The files `cubane.pdb` and `octane.pdb` will be listed.
5. Only the file `octane.pdb` will be listed.

#### *Solution*

4 is the correct answer. `*` matches zero or more characters, so a file name with zero or more characters before a letter `c` and zero or more characters after the letter `c` will be matched.

### Saving to a File in a Loop - Part One

In the `data-shell/molecules` directory, what is the effect of this loop?

```
for alkanes in *.pdb
> do
>   echo $alkanes
>   cat $alkanes > alkanes.pdb
> done
```

1. Prints `cubane.pdb`, `ethane.pdb`, `methane.pdb`, `octane.pdb`, `pentane.pdb` and `propane.pdb`, and the text from `propane.pdb` will be saved to a file called `alkanes.pdb`.

2. Prints `cubane.pdb`, `ethane.pdb`, and `methane.pdb`, and the text from all three files would be concatenated and saved to a file called `alkanes.pdb`.
3. Prints `cubane.pdb`, `ethane.pdb`, `methane.pdb`, `octane.pdb`, and `pentane.pdb`, and the text from `propane.pdb` will be saved to a file called `alkanes.pdb`.
4. None of the above.

*Solution*

1. The text from each file in turn gets written to the `alkanes.pdb` file. However, the file gets overwritten on each loop iteration, so the final content of `alkanes.pdb` is the text from the `propane.pdb` file.

### Saving to a File in a Loop - Part Two

Also in the `data-shell/molecules` directory, what would be the output of the following loop?

```
for datafile in *.pdb
> do
>   cat $datafile >> all.pdb
> done
```

1. All of the text from `cubane.pdb`, `ethane.pdb`, `methane.pdb`, `octane.pdb`, and `pentane.pdb` would be concatenated and saved to a file called `all.pdb`.
2. The text from `ethane.pdb` will be saved to a file called `all.pdb`.
3. All of the text from `cubane.pdb`, `ethane.pdb`, `methane.pdb`, `octane.pdb`, `pentane.pdb` and `propane.pdb` would be concatenated and saved to a file called `all.pdb`.
4. All of the text from `cubane.pdb`, `ethane.pdb`, `methane.pdb`, `octane.pdb`, `pentane.pdb` and `propane.pdb` would be printed to the screen and saved to a file called `all.pdb`.

*Solution*

3 is the correct answer. `>>` appends to a file, rather than overwriting it with the redirected output from a command. Given the output from the `cat` command has been redirected, nothing is printed to the screen.

### 11.11.5 Shell Scripts

#### List Unique Species

Leah has several hundred data files, each of which is formatted like this:

```
2013-11-05,deer,5
2013-11-05,rabbit,22
```

```

2013-11-05,raccoon,7
2013-11-06,rabbit,19
2013-11-06,deer,2
2013-11-06,fox,1
2013-11-07,rabbit,18
2013-11-07,bear,1

```

An example of this type of file is given in `data-shell/data/animal-counts/animals.txt`.

We can use the command `cut -d , -f 2 animals.txt | sort | uniq` to produce the unique species in `animals.txt`. In order to avoid having to type out this series of commands every time, a scientist may choose to write a shell script instead.

Write a shell script called `species.sh` that takes any number of filenames as command-line arguments, and uses a variation of the above command to print a list of the unique species appearing in each of those files separately.

#### *Solution*

```

# Script to find unique species in csv files where species is the second data field
# This script accepts any number of file names as command line arguments

# Loop over all files
for file in $@
do
    echo "Unique species in $file:"
    # Extract species names
    cut -d , -f 2 $file | sort | uniq
done

```

#### Why Record Commands in the History Before Running Them?

If you run the command:

```
$ history | tail -n 5 > recent.sh
```

the last command in the file is the `history` command itself, i.e., the shell has added `history` to the command log before actually running it. In fact, the shell *always* adds commands to the log before running them. Why do you think it does this?

#### *Solution*

If a command causes something to crash or hang, it might be useful to know what that command was, in order to investigate the problem. Were the command only be recorded after running it, we would not have a record of the last command run in the event of a crash.

#### Variables in Shell Scripts



In the `molecules` directory, imagine you have a shell script called `script.sh` containing the following commands:

```
head -n $2 $1
tail -n $3 $1
```

While you are in the `molecules` directory, you type the following command:

```
bash script.sh '*.pdb' 1 1
```

Which of the following outputs would you expect to see?

1. All of the lines between the first and the last lines of each file ending in `.pdb` in the `molecules` directory
2. The first and the last line of each file ending in `.pdb` in the `molecules` directory
3. The first and the last line of each file in the `molecules` directory
4. An error because of the quotes around `*.pdb`

*Solution*

The correct answer is 2.

The special variables `$1`, `$2` and `$3` represent the command line arguments given to the script, such that the commands run are:

```
$ head -n 1 cubane.pdb ethane.pdb octane.pdb pentane.pdb propane.pdb
$ tail -n 1 cubane.pdb ethane.pdb octane.pdb pentane.pdb propane.pdb
```

The shell does not expand `'*.pdb'` because it is enclosed by quote marks. As such, the first argument to the script is `'*.pdb'` which gets expanded within the script by `head` and `tail`.

### Find the Longest File With a Given Extension

Write a shell script called `longest.sh` that takes the name of a directory and a filename extension as its arguments, and prints out the name of the file with the most lines in that directory with that extension. For example:

```
$ bash longest.sh /tmp/data pdb
```

would print the name of the `.pdb` file in `/tmp/data` that has the most lines.

*Solution*

```
# Shell script which takes two arguments:
# 1. a directory name
# 2. a file extension
# and prints the name of the file in that directory
# with the most lines which matches the file extension.

wc -l $1/*. $2 | sort -n | tail -n 2 | head -n 1
```

### Script Reading Comprehension

For this question, consider the `data-shell/molecules` directory once again. This contains a number of `.pdb` files in addition to any other files you may have created. Explain what each of the following three scripts would do when run as `bash script1.sh *.pdb`, `bash script2.sh *.pdb`, and `bash script3.sh *.pdb` respectively.

```
# Script 1
echo *.*

# Script 2
for filename in $1 $2 $3
do
    cat $filename
done

# Script 3
echo $*.pdb
```

#### *Solution*

In each case, the shell expands the wildcard in `*.pdb` before passing the resulting list of file names as arguments to the script.

Script 1 would print out a list of all files containing a dot in their name. The arguments passed to the script are not actually used anywhere in the script.

Script 2 would print the contents of the first 3 files with a `.pdb` file extension. `$1`, `$2`, and `$3` refer to the first, second, and third argument respectively.

Script 3 would print all the arguments to the script (i.e. all the `.pdb` files), followed by `.pdb`. `$*` refers to *all* the arguments given to a shell script.

```
cubane.pdb ethane.pdb methane.pdb octane.pdb pentane.pdb propane.pdb.pdb
```

### 11.11.6 Finding Things

#### Using `grep`

Which command would result in the following output:

and the presence of absence:

1. `grep "of" haiku.txt`
2. `grep -E "of" haiku.txt`
3. `grep -w "of" haiku.txt`
4. `grep -i "of" haiku.txt`

*Solution*

The correct answer is 3, because the `-w` option looks only for whole-word matches. The other options will also match “of” when part of another word.

**Tracking a Species**

Leah has several hundred data files saved in one directory, each of which is formatted like this:

```
2013-11-05,deer,5
2013-11-05,rabbit,22
2013-11-05,raccoon,7
2013-11-06,rabbit,19
2013-11-06,deer,2
```

She wants to write a shell script that takes a species as the first command-line argument and a directory as the second argument. The script should return one file called `species.txt` containing a list of dates and the number of that species seen on each date. For example using the data shown above, `rabbit.txt` would contain:

```
2013-11-05,22
2013-11-06,19
```

Put these commands and pipes in the right order to achieve this:

```
cut -d : -f 2
>
|
grep -w $1 -r $2
|
$1.txt
cut -d , -f 1,3
```

Hint: use `man grep` to look for how to grep text recursively in a directory and `man cut` to select more than one field in a line.

An example of such a file is provided in `data-shell/data/animal-counts/animals.txt`

*Solution*

```
grep -w $1 -r $2 | cut -d : -f 2 | cut -d , -f 1,3 > $1.txt
```

You would call the script above like this:

```
$ bash count-species.sh bear .
```

**Little Women**

You and your friend, having just finished reading *Little Women* by Louisa May Alcott, are in an argument. Of the four sisters in the book, Jo, Meg, Beth, and Amy, your friend thinks that Jo was the most mentioned. You, however, are

certain it was Amy. Luckily, you have a file `LittleWomen.txt` containing the full text of the novel (`data-shell/writing/data/LittleWomen.txt`). Using a `for` loop, how would you tabulate the number of times each of the four sisters is mentioned?

Hint: one solution might employ the commands `grep` and `wc` and a `|`, while another might utilize `grep` options. There is often more than one way to solve a programming task, so a particular solution is usually chosen based on a combination of yielding the correct result, elegance, readability, and speed.

#### *Solution*

```
for sis in Jo Meg Beth Amy
do
    echo $sis:
    grep -ow $sis LittleWomen.txt | wc -l
done
```

Alternative, slightly inferior solution:

```
for sis in Jo Meg Beth Amy
do
    echo $sis:
    grep -ocw $sis LittleWomen.txt
done
```

This solution is inferior because `grep -c` only reports the number of lines matched. The total number of matches reported by this method will be lower if there is more than one match per line.

### Matching and Subtracting

The `-v` option to `grep` inverts pattern matching, so that only lines which do *not* match the pattern are printed. Given that, which of the following commands will find all files in `/data` whose names end in `s.txt` (e.g., `animals.txt` or `planets.txt`), but do *not* contain the word `net`? Once you have thought about your answer, you can test the commands in the `data-shell` directory.

1. `find data -name '*s.txt' | grep -v net`
2. `find data -name *s.txt | grep -v net`
3. `grep -v "temp" $(find data -name '*s.txt')`
4. None of the above.

#### *Solution*

The correct answer is 1. Putting the match expression in quotes prevents the shell expanding it, so it gets passed to the `find` command.

Option 2 is incorrect because the shell expands `*s.txt` instead of passing the wildcard expression to `find`.

Option 3 is incorrect because it searches the contents of the files for lines which do not match “temp”, rather than searching the file names.

### find Pipeline Reading Comprehension

Write a short explanatory comment for the following shell script:

```
wc -l $(find . -name '*.dat') | sort -n
```

*Solution*

1. Find all files with a `.dat` extension recursively from the current directory
2. Count the number of lines each of these files contains
3. Sort the output from step 2. numerically

### Finding Files With Different Properties

The `find` command can be given several other criteria known as “tests” to locate files with specific attributes, such as creation time, size, permissions, or ownership. Use `man find` to explore these, and then write a single command to find all files in or below the current directory that are owned by the user `ahmed` and were modified in the last 24 hours.

Hint 1: you will need to use three tests: `-type`, `-mtime`, and `-user`.

Hint 2: The value for `-mtime` will need to be negative—why?

*Solution*

Assuming that Nelle’s home is our working directory we type:

```
$ find ./ -type f -mtime -1 -user ahmed
```

## 11.12 Key Points

- A shell is a program whose primary purpose is to read commands and run other programs.
- The shell’s main advantages are its high action-to-keystroke ratio, its support for automating repetitive tasks, and its capacity to access networked machines.
- The shell’s main disadvantages are its primarily textual nature and how cryptic its commands and operation can be.
- The file system is responsible for managing information on the disk.
- Information is stored in files, which are stored in directories (folders).
- Directories can also store other directories, which forms a directory tree.

- `cd path` changes the current working directory.
- `ls path` prints a listing of a specific file or directory; `ls` on its own lists the current working directory.
- `pwd` prints the user's current working directory.
- `/` on its own is the root directory of the whole file system.
- A relative path specifies a location starting from the current location.
- An absolute path specifies a location from the root of the file system.
- Directory names in a path are separated with `/` on Unix, but `\\` on Windows.
- `..` means 'the directory above the current one'; `.` on its own means 'the current directory'.
- `cp old new` copies a file.
- `mkdir path` creates a new directory.
- `mv old new` moves (renames) a file or directory.
- `rm path` removes (deletes) a file.
- `*` matches zero or more characters in a filename, so `*.txt` matches all files ending in `.txt`.
- `?` matches any single character in a filename, so `? .txt` matches `a.txt` but not `any.txt`.
- Use of the Control key may be described in many ways, including `Ctrl-X`, `Control-X`, and `^X`.
- The shell does not have a trash bin: once something is deleted, it's really gone.
- Most files' names are `something.extension`. The extension isn't required, and doesn't guarantee anything, but is normally used to indicate the type of data in the file.
- Depending on the type of work you do, you may need a more powerful text editor than Nano.
- `cat` displays the contents of its inputs.
- `head` displays the first 10 lines of its input.
- `tail` displays the last 10 lines of its input.
- `sort` sorts its inputs.
- `wc` counts lines, words, and characters in its inputs.

- `command > file` redirects a command's output to a file (overwriting any existing content).
- `command >> file` appends a command's output to a file.
- `<` operator redirects input to a command
- `first | second` is a pipeline: the output of the first command is used as the input to the second.
- The best way to use the shell is to use pipes to combine simple single-purpose programs (filters).
- A `for` loop repeats commands once for every thing in a list.
- Every `for` loop needs a variable to refer to the thing it is currently operating on.
- Use `$name` to expand a variable (i.e., get its value). `${name}` can also be used.
- Do not use spaces, quotes, or wildcard characters such as `'*or?'` in file-names, as it complicates variable expansion.
- Give files consistent names that are easy to match with wildcard patterns to make it easy to select them for looping.
- Use the up-arrow key to scroll up through previous commands to edit and repeat them.
- Use `Ctrl-R` to search through the previously entered commands.
- Use `history` to display recent commands, and `!number` to repeat a command by number.
- Save commands in files (usually called shell scripts) for re-use.
- `bash filename` runs the commands saved in a file.
- `$@` refers to all of a shell script's command-line arguments.
- `$1`, `$2`, etc., refer to the first command-line argument, the second command-line argument, etc.
- Place variables in quotes if the values might have spaces in them.
- Letting users decide what files to process is more flexible and more consistent with built-in Unix commands.
- `find` finds files with specific properties that match patterns.
- `grep` selects lines in files that match patterns.
- `--help` is an option supported by many bash commands, and programs that can be run from within Bash, to display more information on how to use these commands or programs.

- `man` command displays the manual page for a given command.
- `$(command)` inserts a command's output in place.







## Chapter 12

# A Branching Workflow

### 12.1 Questions

- How can a growing number of people coordinate work on a single project?

### 12.2 Objectives

- Explain what rebasing is and use it interactively to collapse a sequence of commits into a single commit.
- Describe a branch-per-feature workflow and explain why to use it.
- Describe what a repository tag is and create an annotated tag in a Git repository.

### 12.3 Introduction

A common Git workflow for single-author/single-user projects is:

1. Make some changes, running `git add` and `git commit` as you go along.
2. `git push` to send those changes to a remote repository for backup and sharing.
3. Switch computers (e.g., go from your laptop to the cluster where you're going to process full-sized data sets).
4. `git pull` to get your changes from the remote repository, and possibly resolve merge conflicts if you made some changes on the cluster the last time you were there but forgot to push them to the remote repo.

Occasionally, you'll run `git checkout -- .` or `git reset --hard VERSION` to discard changes. Less frequently, you'll use `git diff -r abcd1234 path/to/file` or `git show abcd1234:path/to/file` to look at what's changed in files, or `git checkout abcd1234 path/to/file` followed by `git add` and `git commit` to restore an old version of a file. And every once in a while, you'll jenny the repository and clone a new copy because everything is messed up.

This workflow essentially uses Git to do backups and as a replacement for `scp` (or `ftp`, if you're of an older generation). But sometimes you have to work on several things at once in a single project, or need to set aside current work for a high-priority interrupt. And even if those situations don't arise, this workflow doesn't provide guidance for collaborating with others.

This lesson introduces a few working practices that will help you do these things. Along the way, it shows you how to use some of Git's more advanced capabilities. As with everything involving Git, the syntax for commands is more complicated than it needs to be (Perez De Rosso and Jackson (2013), Perez De Rosso and Jackson (2016)), but if used carefully, they can make you much more productive.

## 12.4 How can I use branches to manage development of new features?

Experienced developers tend to use a branch-per-feature workflow, even when they are working on their own. Whenever they want to create a new feature, they create a new branch from `master` and do the work there. When it's done, they merge that branch back into 'master'.

More specifically, a branch-per-feature workflow looks like this:

1. `git checkout master` to make sure you're working in the `master` branch.
2. `git pull origin master` (or `git pull upstream master`) to make sure that your starting point is up to date (i.e., that the changes you committed while working on the cluster yesterday are now in your local copy, or that the fix your colleague did last week has been incorporated).
3. `git checkout -b name-of-feature` to create a new branch. The branch's name should be as descriptive as the name of a function or of a file containing source code: "`stuff`" or "`fixes`" saves you a couple of seconds of typing when you create the branch, and will cost you much more time than that later on when you're juggling half a dozen branches and trying to figure out what's in each.
4. Do the work you want to do. If some other task occurs to you while you're doing it—for example, if you're writing a new function and realize that the documentation for some other function should be updated—*don't* do

## 12.5. HOW CAN I SWITCH BETWEEN BRANCHES WHEN WORK IS ONLY PARTLY DONE?229

that work in this branch just because you happen to be there. Instead, commit your changes, switch back to **master**, create a new branch for the other work, and carry on.

5. Periodically, you may `git push origin name-of-feature` to save your work to GitHub.
6. When the work is done, you `git pull origin master` to get any changes that anyone else has put in the master copy into your branch and merge any conflicts. This is an important step: you want to do the merge *and test that everything still works* in *your* branch, not in **master**.
7. Finally, use `git checkout master` to switch back to the **master** branch on your machine and `git merge name-of-feature master` to merge your changes into **master**. Test it one more time, and then `git push origin master` to send your changes to the remote repository for everyone else to admire and use.

It takes longer to describe this workflow than to use it. Right now, for example, the repository for this material currently has six active branches: **master** (which is the published version), **explain-branch-per-feature** (where this is being written), and four others created from **master** as reminders to fix a glossary entry, rewrite the description in `README.md` of how to generate a PDF of these notes, and so on. This approach essentially uses branches as a to-do list, which takes no more time than typing a brief note about the feature into a text file or writing it in a lab notebook, and is much easier to track. Chapter 19 will look at better ways to manage this.

## 12.5 How can I switch between branches when work is only partly done?

So what happens if we are in the middle of making changes on branch X and realize that we need to do some reorganizing that has nothing to do with the feature we're currently working on? The easy answer is to commit our half-finished work, switch branches, do the other task, and then switch back. However, that workflow leaves a lot of commits in our history corresponding to half-done work, and ideally we want every commit to move the system from one useful state to another. And if we make a dozen commits while working on a feature and then merge that branch into **master**, the overall history of the project can become very difficult to read.

There is a command called `git stash` that will temporarily save work on a branch. We don't recommend using it, since it creates yet another place where valuable information might be stored (or forgotten). Instead, commit changes in the branch as you work, then squash those commits into one single comprehensible commit using the command `git rebase` (discussed in the next sec-

tion). After squashing, you can merge that single large (meaningful) commit into **master**.

Like many things involving Git, this is easier to understand with a diagram. Suppose we have created a branch to work on a feature:

FIXME: figure

We make several small changes, then realize we need to fix something else, so we commit those changes:

FIXME: figure

switch to **master**, and create a new branch to do the other work. This happens several times, eventually leaving the repository in this state:

FIXME: figure

If we were to merge now, four commits would wind up in the history of the **master** branch. Instead, we squash those four commits into one:

FIXME: figure

and then merge that:

FIXME: figure

## 12.6 How can I keep my project's history clean when working on many branches?

The workflow described above depends on rebasing, which means moving or combining some commits from one branch to another. **git rebase** is a powerful command: it can replay changes made in one branch on top of changes made to another:

FIXME: figure

or collapse several consecutive commits into a single commit, as we saw in the previous section. We will only use it for the latter ability.

The command we want is **git rebase -i BASE**, where **BASE** is the most recent commit *before* the sequence to be compressed (i.e., the one that everything else will be based on). Many people find this confusing, and frequently ask **git rebase** to start with the first commit they want changed rather than the last one that they don't.

When we run **git rebase -i**, it brings up a display of recent commits in the same editor we would use for writing commit messages:

FIXME: figure

## 12.7. HOW CAN I MAKE IT EASY FOR PEOPLE TO REVIEW MY WORK BEFORE I MERGE IT?231

The help text in this display tells us what we can do; in most cases, people choose to **pick** the first commit and **squash** the rest. When we save this file, Git immediately launches the editor again to show us the combination of recent changes and messages:

FIXME: figure

We can now edit this text to create a commit message for the unified commit. After we have done this, `git log` will only show us this single commit, because it's the only one left in our history:

FIXME: figure

`git rebase` is a complex command, and if we have merged changes from other branches into the branch we're rebasing, Git can become confused (or rather, it can confuse us). A good rule to follow is, "Don't rebase branches that are shared with other people." Instead, do all the work required for a feature, rebase, and then merge in recent changes from **master** and merge back to **master**:

FIXME: figure

What if you have pushed a branch to GitHub (or elsewhere) and then combine the commits to change its history as shown below:

FIXME: figure

In this case, Git realizes that your merge would lose information and prevents the operation from going through. You can use `git push --force` to overwrite the remote history, but this is usually a sign that you should have done something differently a while back. Remember, `git push --force` will also overwrite any work that other people have pushed to the repository, so it's a good way to end friendships.

## 12.7 How can I make it easy for people to review my work before I merge it?

The biggest benefit of having a second person work on a programming project is not getting twice as much code written, but having a second pair of eyes look at the software. Chapter ?? will discuss how to do code review, but a good first step depends on using a branch-per-feature workflow.

In order for someone to review your change, that change and the original have to somehow be put side-by-side. The usual way to do this is to create a pull request, which is a marker saying, "Someone would like to merge branch A into branch X". The pull request does *not* contain the changes: instead, it points at two branches, so that if either branch is changed, the differences displayed are always up to date.

FIXME: figure

Pull requests can be created between two branches in a single repository, or between branches of different repositories. The latter is common in projects with many contributors: each contributor works in a branch in their own fork of the repository, and when they're done, they create a pull request offering to merge their changes into the **master** branch of the original repository. In order to update their own repository, they pull changes from the original:

FIXME: figure

A pull request can store more than the location of the source and destination branch. In particular, it can store comments that people have made about the proposed merge. GitHub and other forges allow users to add comments on the pull request as a whole, or on particular lines, and can mark old comments as out of date if the lines of code the comment is attached to are updated.

Pull requests aren't just for code. If you are using a typesetting language like Markdown or LaTeX to write your papers, you can create pull requests for changes so that your colleagues can comment on them, then revise and push again when you incorporate their suggestions. This workflow has all the benefits of adding comments to Google Docs, but scales to large projects or large numbers of contributors. The downside, of course, is that you have to use a typesetting language rather than a WYSIWYG editor, because the programmers who created Git and other version control systems still don't support anything that couldn't be put on a punchcard in 1965.

## 12.8 What *is* a feature?

But what is a “feature”, exactly? What's large enough to merit creation of a new branch? These rules make sense for small projects with or without collaborators:

1. Anything cosmetic that is only one or two lines long can be done in **master** and committed right away. “Cosmetic” means changes to comments or documentation: nothing that affects how code runs—not even a simple variable renaming—is done this way, because experience has taught that things that aren't supposed to often do.
2. A pure addition that doesn't change anything else is a feature and goes into a branch. For example, if you run a new analysis and save the results, that should be done on its own branch because it might take several tries to get the analysis to run, and you might interrupt yourself several times to fix things that you've discovered aren't working.
3. Every change to code that someone might want to undo later in one step gets done as a feature. For example, if a parameter is added to a function, every call to the function has to be updated; since neither alteration makes



sense without the other, it's considered a single feature and should be done in one branch.

The hardest thing about using a branch-per-feature workflow is actually doing it for small changes. As the first point in the list above suggests, most people are pragmatic about this on small projects; on large ones, where dozens of people might be committing, even the smallest and most innocuous change needs to be in its own branch so that it can be reviewed (which we discuss below).

The other thing that's hard to do with a branch-per-feature workflow is a major code reorganization. If many files are being moved, renamed, and altered in order to restructure the project, merging branches where those changes *haven't* been made can be tedious and error-prone. The solution is to not get into this situation: as Chapter ?? says, code should be reorganized in many small steps, not one big one.

## 12.9 How can I label specific versions of my work?

A tag is a permanent label on a particular state of the repository. Tags are theoretically redundant, since the commit hash identifies that state as well, but commit hashes are (deliberately) random and therefore hard to remember or find.

Experienced developers Use annotated tags to mark every major event in the project's history. These tags are called "annotated" because they allow their creator to specify a message, just like a commit. Research projects often use **report-date-event** for tag names, such as **jse-2018-06-23-response** or **pediatrics-2018-08-15-summary**. If you do this, please don't tempt fate by calling something **-final**.

Most software projects use semantic versioning for software releases, which produces three-part version numbers **major.minor.patch**:

- Increment **major** every time there's an incompatible externally-visible change.
- Increment **minor** when adding functionality without breaking any existing code.
- Increment **patch** for bug fixes that don't add any new features ("now works as previously documented").

Simple projects only tag the **master** branch because everything that is finished is merged to **master**. Larger software projects may create a branch for each released version and do minor or patch updates on that branch, but this outside the scope of this lesson.

## 12.10 Summary

FIXME: create concept map for workflow

## 12.11 Exercises

FIXME: exercises for branching

## 12.12 Key Points

- Create a new branch for every feature, and only work on that feature in that branch.
- Always create branches from **master**, and always merge to **master**.
- Use rebasing to combine several related commits into a single commit before merging.

## Chapter 13

# Configuring Software

### 13.1 Questions

- How can I make it easy for users to configure software?

### 13.2 Objectives

- Describe the four levels of configuration typically used by robust software.
- Explain what an overlay configuration is.
- Explain why nested configuration options are usually not a good idea.
- Add flat overlay configuration to a small application.

### 13.3 Introduction

A program that does exactly the same thing every time we run it isn't as useful as one that can work on different files or analyze data with different thresholds. Software of all kinds needs to be controlled; some things change more often than others, so we need a simple, uniform way to specify some options and leave others alone.

The modern Unix convention is to provide four levels of configuration:

1. A system-wide configuration file for general settings.
2. A user-specific configuration file for personal preferences.
3. A job-specific file with settings for a specific run.
4. Command-line options to change things that commonly change.

This is sometimes called overlay configuration because each level overrides the ones above it: the user’s configuration file overrides the system settings, the job configuration overrides the user’s defaults, and the command-line options overrides that.

## 13.4 How can I handle command-line flags consistently?

Modern Python programs use the `argparse` library for handling command-line arguments, but the older and simpler `getopt` library will illustrate the core ideas, so we will use it.

`getopt` works by matching a specification of what flags are allowed against a list of actual command-line arguments. In simplest form, the spec is string listing all the single-letter flags, with colons showing the ones that take arguments. For example, the string `b:q` means “the `-b` flag takes an argument and the `-q` flag doesn’t”.

The list of actual arguments is almost always `sys.argv[1:]`, i.e., all of the command-line arguments except for the name of the program itself. (The name `argv` stands for “argument vector”, and is a holdover from the days of C.) Given the spec and the list of actual arguments, `getopt` returns two lists: the first is the (flag, argument) pairs it matched, and the second is everything else—typically a list of files to be processed.

```
from getopt import getopt

args = ['-q', '-b', '/tmp/log.txt', 'file1.txt', 'file2.txt']
options, extras = getopt(args, 'b:q')
print('options is', options)
print('extras is', extras)
```

```
options is [('-q', ''), ('-b', '/tmp/log.txt')]
extras is ['file1.txt', 'file2.txt']
```

Once we have these two lists, we can define default values for our configuration and then loop over the (flag, argument) pairs to override those values:

```
# Defaults.
logfile = None
quiet = False

# Override based on command-line options.
for (opt, arg) in options:
    if opt == '-b':
        logfile = arg
```

```

elif opt == '-q':
    quiet = True
else:
    assert False, 'unrecognized option {}'.format(opt)
print('Log file is {} and quiet is {}'.format(logfile, quiet))

```

Log file is /tmp/log.txt and quiet is True

This pattern is called set and override, and makes programs easier to understand by putting all of the default settings in one place. Since we will often want to pass those settings into functions that do the actual work, it's very common to put the entire configuration in a dictionary so that we have a single configuration object to pass around. It's also common to check that some configuration values aren't accidentally being set twice. After rearranging our code a little, we get this:

```

import sys
from getopt import getopt

# Defaults.
settings = {
    'logfile' : None,
    'quiet' : False
}

# Override based on command-line options.
options, extras = getopt(sys.argv[1:], 'b:q')
for (opt, arg) in options:
    if opt == '-b':
        assert settings['logfile'] is None, 'cannot set logfile twice'
        settings['logfile'] = arg
    elif opt == '-q':
        settings['quiet'] = True
    else:
        assert False, 'unrecognized option {}'.format(opt)

# Display.
print('Log file {}'.format(settings['logfile']))
print('Quiet {}'.format(settings['quiet']))
print('Extras {}'.format(extras))

```

which we can run like this:

```
$ python getopt_dict.py -q first.txt
```

```

Log file None
Quiet True
Extras ['first.txt']

```

We really shouldn't use `assert` to handle errors here; Chapter E.3.5 will explore a better approach.

## 13.5 What do I do when I run out of memorable single-letter flags?

Taschuk's Third Rule says, "Make common operations easy to control." Taschuk and Wilson (2017) so that users can control everything from a shell script without having to create temporary configuration files. However, there are only so many single-letter flags available, and using `-b` to specify the name of a log file is hardly intuitive.

To allow for this, `getopt` handles another kind of flag: a double dash followed by a longer name like `--erase-temp-files`. These flags can take a single argument like their single-letter siblings. To tell `getopt` what long names it should recognize, we give the function an extra list and use `=` as a suffix to indicate if the option takes an argument.

```
from getopt import getopt

args = ['-q', '--logfile', '/tmp/log.txt', '--overwrite', 'file1.txt', 'file2.txt']
options, extras = getopt(args, 'b:q', ['logfile=', 'overwrite'])
print('options is', options)
print('extras is', extras)

options is [('-q', ''), ('--logfile', '/tmp/log.txt'), ('--overwrite', '')]
extras is ['file1.txt', 'file2.txt']
```

It's common to use single-letter flags for the most frequently changed options and long names for things that are changed less frequently, and to provide long-name aliases for the single-letter flags (e.g., to have `--all-files` mean the same thing as `-a`).

## 13.6 How can I manage configuration files consistently?

Controlling programs from the command line is useful, but complex programs can have many different configuration options, and it's very useful to be able to save settings in a file for later reference (and reproducibility). Enabling a program to read its configuration from a file also allows users to set values once and then not worry about them, which is particularly useful when they're installing the software on their own computer and want to put temporary files

### 13.6. HOW CAN I MANAGE CONFIGURATION FILES CONSISTENTLY?239

in a different location or change the value of the alpha parameter for fitting curves.

Programmers have invented many formats for configuration files, so please do not create your own. One possibility is to write the configuration as a Python data structure and then load it as if it was a library. This is clever, but it's hard for tools in other languages to process. Programers are also fond of JSON, which is a subset of the syntax that JavaScript uses for data structures, but that involves a lot of curly braces. A third option is the Windows INI format, which is laid out like this:

```
[section_1]
key_1=value_1
key_2=value_2

[section_2]
key_3=value_3
key_4=value_4
```

INI files are simple to read and write, but the format is slowly falling out of use. What seems to be replacing it is YAML, which stands for “Yet Another Markup Language”. Since YAML is used in GitHub Pages, and (unlike JSON) allows comments, we'll explore it in this section.

Here's a sample configuration file:

```
# Example configuration file
logfile: "/tmp/log.txt"
quiet: false
overwrite: false
fonts:
- Verdana
- Serif
```

And here's a short Python program that reads and prints that configuration:

```
import yaml

with open('config.yaml', 'r') as reader:
    config = yaml.load(reader)
print(config)
```

```
{'logfile': '/tmp/log.txt', 'quiet': False, 'overwrite': False, 'fonts': ['Verdana', 'Serif']}
```

Simple YAML files are simple to write:

1. Lines starting with # are comments.
2. A line **key: value** defines a value for the given key. Values can be numbers, **true** or **false**, or quoted strings. (Strings don't actually have to be quoted in every case, but the file is a lot easier to understand if you

always use quotes.)

3. A point-form list underneath a key becomes an array of values.

When a file like this is read in Python, the result is a dictionary. YAML allows nested keys and lists, but if you need them, you're probably doing something wrong Xu et al. (2015): most users never use most configuration options and find their presence confusing.

## 13.7 How can I implement overlay configuration?

We said at the start that programs often have system-wide, per-user, and per-job configuration files, with each overriding values from the one(s) before and command-line parameters overriding the rest. We can implement this using `dict.update`, which updates one dictionary with values from another:

```
def get_full_configuration(filenamees, command_line={}):
    """
    Overlay configuration files and command-line parameters,
    returning configuration object.
    """
    result = {}
    for f in filenamees:
        with open(f, 'r') as reader:
            config = yaml.load(reader)
            result.update(config)
    result.update(command_line)
    return result
```

This function creates an empty dictionary to hold settings. It then reads each specified configuration file in turn and updates the result dictionary with whatever it found in that file. If a file defines values that were previously defined in an earlier file, the `update` method call automatically overwrites the older values. We end by overriding what we read from the files with whatever was given on the command line; we'll have to convert `getopt`'s output to a dictionary, but that's straightforward—the only trick is to match the command-line flag (like `-q`) to the configuration variable name (like `quiet`):

```
def getopt_to_dict(pairs, names):
    """
    Convert [(flag, value)...] pairs and {flag: config...} names
    to dictionary.
    """
    result = {}
    for (key, value) in pairs:
```



```

        result[names[key]] = value
    return result

```

We can test this with these three configuration files (we have lined up corresponding values to make them easier to see):

system.yml

user.yml

job.yml

quiet: true

quiet: false

logfile: “/tmp/log.txt”

logfile: “./complaints.txt”

using this test program:

```

import sys
from getopt import getopt
from util import get_full_configuration, getopt_to_dict

config_files = ['system.yml', 'user.yml', 'job.yml']

options, extras = getopt('b:q', sys.argv[1:])
options = getopt_to_dict(options, {'-b': 'logfile', '-q': 'quiet'})
config = get_full_configuration(config_files, options)
print(config)

```

and this command line:

```
$ python test-config.py -q
```

```
{'quiet': False, 'logfile': './complaints.txt'}
```

## 13.8 How can I find configuration files?

Our configuration files will usually not all be in the same directory. System-wide settings for an application called `app` are often stored in `/etc/app.yml`. Alternatively, some programs will set an environment variable `APP` to the name of the installation directory, and then read the system configuration file from `$APP/app.yml`. We can use `os.getenv('APP')` to get the value of the environment variable `APP`, then append `app.yml` and load that. (Older programs often use the name `app.rc`, where “rc” stands for “resource control”.)

Similarly, we can get personal settings from `$HOME/.app.yml`; the leading `.` hides the configuration file from `ls`. Finally, per-job settings can come from `app.yml` in the current directory, where again “app” is replaced with the name of the program. Here’s a utility routine that constructs and checks these filenames:

```
def find_config_files(name):  
    '''  
    Construct a list of configuration files for the named application.  
    '''  
    app_yaml = name + '.yaml'  
    locations = [(name.upper(), app_yaml),  
                 ('HOME', '.' + app_yaml),  
                 ('PWD', app_yaml)]  
    result = []  
    for (var, filename) in locations:  
        value = os.getenv(var)  
        if value:  
            path = os.path.join(value, filename)  
            if os.path.isfile(path):  
                result.append(path)  
    return result
```

Note that we use `os.path.isfile` to check that file exists before trying to read it.

### 13.9 How can I keep a record of the actual configuration that produced particular results?

Careful record keeping is essential to reproducible science, and if *we* are careful, the computer can do the record keeping. We can save the entire (merged) configuration object for a particular run of a program using `yaml.dump`. If we have written our configuration functions correctly, this will let us re-create configuration on another machine even if it has different default settings. The test is whether our program can load a dumped configuration, then dump it again and get the same result.

If we’re going to do this, we should always include a version number as a field in the dumped configuration; our program should also print this out when given a `--version` flag. We need this because how we interpret options will change over time, and if you don’t know what the version of the program was, we’ll have to guess what options mean.

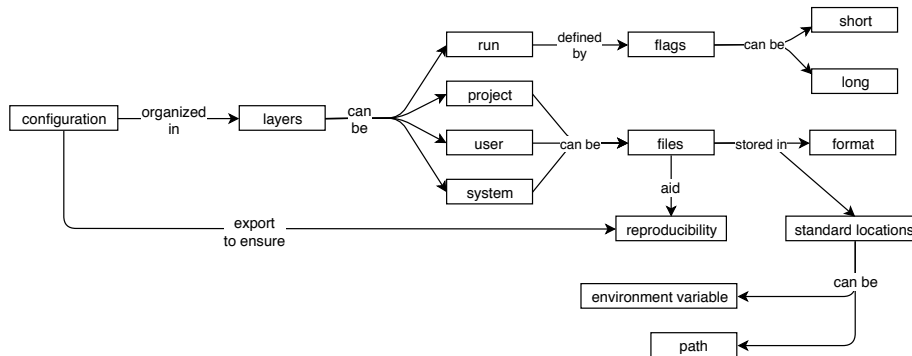


Figure 13.1: Configuration Concept Map

## 13.10 Summary

- Many tool also allow longer command-line options to control everything to support scripting
  - Out of scope of this lesson
- Every tool acts as if it was the only extra thing your project needed
  - So you wind up with lots of configuration files littering your root directory
  - Tempting to move them all to `./etc/`, but tools don't know to look there

## 13.11 Exercises

### 13.11.1 Dump and check configuration

- FIXME: how to handle lack of order in keys? (answer: use deep equality)

### 13.11.2 Read configuration solely from file

- FIXME: implement flag to read entire configuration from file *and nothing else*.

### 13.11.3 Implement a usage function

- FIXME: implement `usage` function to display message and exit instead of using `assert`.

#### 13.11.4 Read standard input if no files specified

- FIXME: how to handle reading from stdin/stdout if no files specified.

### 13.12 Key Points

- Use short command-line parameters to set commonly-changed options.
- Allow long command-line parameters to change all other options to facilitate scripting.
- Use the `getopt` library to parse command-line flags.
- Read a system-wide configuration file, then a user configuration file, then a job configuration file.
- Format configuration files using YAML.
- Dump all configuration values in the same format used for input on request.
- Include the software version number in the dumped configuration.

## Chapter 14

# Automating Analyses

### 14.1 Questions

- How can I make my analyses easier to reproduce?
- How can I make it easier to repeat analyses when I get new data, or when my data or scripts change?

### 14.2 Objectives

- Explain what a build tool is and how build tools aid reproducible research.
- Describe and identify the three parts of a Make rule.
- Write a Makefile that re-runs a multi-stage data analysis.
- Explain and trace how Make chooses an order in which to execute rules.
- Explain what phony targets are and define a phony target.
- Explain what automatic variables are and correctly identify three commonly-used automatic variables.
- Rewrite Make rules to use automatic variables.
- Explain why and how to write a pattern rule in a Makefile.
- Rewrite Make rules to use patterns.
- Define variables in a Makefile explicitly and by using functions.
- Make a self-documenting Makefile.

### 14.3 Introduction

As Section 10.4 said, Zipf's Law states that the second most common word in a body of text appears half as often as the most common, the third most common

appears a third as often, and so on. The analyses we want to do include:

- Analyze one input file to see how well it conforms to Zipf's Law.
- Analyze multiple input files to how well then conform in aggregate.
- Plot individual and aggregate word frequency distributions and expected values.

The project we have inherited as a starting point contains the following:

1. The books we are analyzing are in `data/title.txt`.
2. A program called `bin/countwords.py` can read a text file and produce a CSV file with two columns: the words in the text and the number of times that word occurs.
3. `bin/countwords.py` can analyze several files at once if we provide many filenames on the command line or concatenate the files and send them to standard input in a pipeline using something like `cat data/*.txt | bin/countwords.py`.
4. Another program, `bin/plotcounts.py`, will create a visualization for us that shows word rank on the X axis and word counts on the Y axis. (It doesn't show the actual words.)
5. A third program, `bin/collate.py`, takes one or more CSV files as input and merges them by combining the counts for words they have in common.
6. Finally, `bin/testfit.py` will compare actual distributions against expectations and give a fitting score.

It's easy enough to run these programs by hand if we only want to analyze a handful of files, but doing this becomes tedious and error-prone as the number of files grows. Instead, we can write a shell script or another Python program to do multi-file analyses. Doing this documents the pipeline so that our colleagues (and future selves) can understand it, and enables us to re-do the entire analysis with a single command if we get new data or change our methods, parameters, or libraries. It also prevents us from making lots of little errors: there's no guarantee we'll get the script right the first time, but once we've fixed it, it will stay fixed.

However, re-running the entire analysis every time we get a new file is inefficient: we don't need to re-count the words in our first thousand books when we add the thousand and first. This isn't a problem when calculations can be done quickly, but many can't, and anyway, the point of this chapter is to introduce a tool for handling cases in which we really want to avoid doing unnecessary work.

What we want is a way to describe which files depend on which other files and how to generate or update a file when necessary. This is the job of a build tool. As the name suggests, a build tool's job is to build new files out of old ones. The most widely used build tool, Make, was written in 1976 to recompile programs (which at the time was a slow process). GNU Make is a free, fast, and well-documented version of Make; we will use it throughout this book.

### Alternatives

Many better build tools have been developed since Make—so many, in fact, that none has been able to capture more than a small fraction of potential users. Snakemake has a lot of fans, and a future version of this tutorial might well use it.

This introduction based on the Software Carpentry lesson on Make maintained by Gerard Capes and on Jonathan Dursi's introduction to pattern rules.

## 14.4 How can I update a file when its prerequisites change?

Make is based on three key ideas:

1. The operating system automatically records a timestamp every time a file is changed. By checking this, Make can tell whether files are newer or older than other files.
2. A programmer writes a Makefile to tell Make how files depend on each other. For example, the Makefile could say that `results/moby_dick.csv` depends on `data/moby_dick.txt`, or that `plots/aggregate.svg` depends on all of the CSV files in the `results/` directory.
3. The Makefile includes shell commands to create or update files that are out of date. For example, it could include a command to (re-)run `bin/countwords.py` to create `results/moby_dick.csv` from `data/moby_dick.txt`. (Make's use on shell commands is one reason for its longevity, since it allows programmers to write tools for updating files in whatever language they want.)

Let's start by creating a file called `Makefile` that contains the following three lines:

```
# Regenerate results for "Moby Dick"
results/moby_dick.csv : data/moby_dick.txt
    python bin/countwords.py data/moby_dick.txt > results/moby_dick.csv
```

The `#` character starts a comment, which runs to the end of a line (just as it does in Python or R). `results/moby_dick.csv` is the target of a rule, i.e., something that may need to be created or updated. Every rule in a Makefile has one or more targets, and must be written flush against the left margin.

`data/moby_dick.txt` is a prerequisite in that rule, i.e., something that the target of the rule depends on. A single colon separates the target from its prerequisites, and a rule can have any number of prerequisites—we'll see examples soon.

The indented line that uses Python to run `bin/countwords.py` is the rule's

action. It creates or updates the target when it is out of date. A rule can have any number of actions, but they *must* be indented by a single tab character. Notice that the output of `bin/countwords.py` is redirected using `>` to create the output file: we will look at modifying the script in Chapter 13 so that it can take the name of an output file as an argument.

Together, the three parts of this rule tell Make when and how to re-create `results/moby_dick.csv`. To test that it works, run this command in the shell:

```
$ make
```

Make automatically looks for a file called `Makefile` and checks the rules it contains. In this case, one of three things will happen:

1. Make won't be able to find the file `data/moby_dick.csv`, so it will run the script to create it.
2. Make will see that `data/moby_dick.txt` is newer than `results/moby_dick.csv`, in which case it will run the script to update the results file.
3. Make will see that `results/moby_dick.csv` is newer than its prerequisite, so it won't do anything.

In the first two cases, Make will show the command it runs, along with anything the command prints to the screen via standard output or standard error. In this case, there is no screen output, so we only see the command.

### Indentation Errors

If `Makefile` contains spaces instead of tabs to indent the rule's action, we will see an error message like this:

```
Makefile:3: *** missing separator. Stop.
```

The requirement to use tabs is a legacy of Make's origins as a student intern project, and no, I'm not kidding.

If we run `make` again right away it doesn't re-run our script because we're in situation #3 from the list above: the target is newer than its prerequisites, so no action is required. We can check this by listing the files with their timestamps, ordered by how recently they have been updated:

```
$ ls -l -t data/moby_dick.txt results/moby_dick.csv
-rw-r--r--  1 gvwilson  staff   219107 31 Dec 08:58 results/moby_dick.csv
-rw-r--r--  1 gvwilson  staff  1276201 31 Dec 08:58 data/moby_dick.txt
```

When Make sees that a target is newer than its prerequisites it displays a message like this:

```
make: `results/moby_dick.csv' is up to date.
```

To test that Make is actually doing the right thing, we can:

1. Delete `results/moby_dick.csv` and type `make` again (situation #1).



2. Run the command `touch data/moby_dick.txt` to update the timestamp on the source file, then run `make` (situation #2).

## 14.5 How can I tell Make where to find rules?

We don't have to call our file of rules `Makefile`. If we want, we can rename the file `single_rule.mk` and then run it with `make -f single_rule.mk`. Most people don't do this in real projects, but in a lesson like this, which includes many example Makefiles, it comes in handy.

Using `-f` doesn't change our working directory. If, for example, we are in `/home/gvwilson/still-magic` and run `make -f src/automate/single_rule.mk`, our working directory remains `/home/gvwilson/still-magic`. This means that Make will look for the rule's prerequisite in `/home/gvwilson/still-magic/data/moby_dick.txt`, not in `/home/gvwilson/still-magic/src/automate/data/moby_dick.txt`.

FIXME: figure

## 14.6 How can I update multiple files when their prerequisites change?

Our Makefile isn't particularly exciting so far. Let's add another rule to the end and save the result as `double_rule.mk`:

```
# Regenerate results for "Moby Dick"
results/moby_dick.csv : data/moby_dick.txt
    python bin/countwords.py data/moby_dick.txt > results/moby_dick.csv

# Regenerate results for "Jane Eyre"
results/jane_eyre.csv : data/jane_eyre.txt
    python bin/countwords.py data/jane_eyre.txt > results/jane_eyre.csv
```

When we ask Make to run this file:

```
$ make -f double_rule.mk
```

we get this rather disappointing message:

```
make: `results/moby_dick.csv' is up to date.
```

Nothing happens because by default Make only attempts to update the first target it finds in the Makefile, which is called the default target. To update something else, we need to tell Make we want it:

```
$ make -f double_rule.mk results/jane_eyre.csv
```

This time Make runs:

```
python bin/countwords.py data/jane_eyre.txt > results/jane_eyre.csv
```

## 14.7 How can I get rid of temporary files that I don't need?

A phony target in a Makefile is one that doesn't correspond to any files and doesn't have any prerequisites. Phony targets are just a way to save useful commands in a Makefile, but saying “just” is a bit misleading: what we're actually doing is recording all of the steps in our workflow, even if those steps don't create or update files.

For example, let's add another target to our Makefile to delete all of the files we have generated. By convention this target is called `clean`, and ours looks like this:

```
# Remove all generated files.
clean :
    rm -f results/*.csv
```

(The `-f` flag to `rm` means “force removal”. When we use it, `rm` won't complain if the files it's trying to remove are already gone.) Let's run Make:

```
$ make -f pipeline.mk clean
```

and then use `ls` to list the contents of `results`. Sure enough, it's empty.

Phony targets are useful as a way of documenting actions in a project, but there's a catch. Use `mkdir` to create a directory called `clean`, then run `make -f clean.mk clean`. Make will print:

```
make: `clean' is up to date.
```

The problem is that Make finds something called `clean` and assumes that's what the rule is referring to. Since the rule has no prerequisites, it can't be out of date, so no actions are executed.

There are two ways to solve this problem. The first is to make sure we don't have phony targets with the same names as files or directories. That works as long as our project is small and we're paying attention, but as the project grows, or we're rushing to meet a deadline or have inherited the project from someone else and don't realize that this might be a problem, it's bound to fail at exactly the worst time.

A much better solution is to tell Make that the target is phony by putting this in the Makefile:

```
.PHONY : clean
```

Most people declare all of their phony targets together near the top of the file, though some put the `.PHONY` declarations right before the rules they refer to. As with most other rules about programming style (Chapter ??), consistency matters more than exactly what you do.

## 14.8 How can I make a target depend on several prerequisites?

Right now, our Makefile says that each result file depends only on the corresponding data file. That's not accurate: in reality, each result also depends on the script used to generate it. If we change our script, we ought to regenerate our results and then check to see if they've changed. (We can rely on version control to tell us that.)

Here's a modified version of the Makefile in which each result depends on both the data file and the script:

```
.PHONY: clean

# Regenerate results for "Moby Dick"
results/moby_dick.csv : data/moby_dick.txt bin/countwords.py
    python bin/countwords.py data/moby_dick.txt > results/moby_dick.csv

# Regenerate results for "Jane Eyre"
results/jane_eyre.csv : data/jane_eyre.txt bin/countwords.py
    python bin/countwords.py data/jane_eyre.txt > results/jane_eyre.csv

# ...clean target as before...
```

We can test this by touching the script and then making one or the other result:

```
$ touch bin/countwords.py
$ make -f depend_on_script.mk results/jane_eyre.csv

python bin/countwords.py data/jane_eyre.txt > results/jane_eyre.csv
```

## 14.9 How can I reduce the amount of typing I have to do?

The name of our script now appears four times in our Makefile, which will make for a lot of typing if and when we decide to move it or rename it. We can fix that by defining a variable at the top of our file to refer to the script, then using that variable in our rules:

```
.PHONY: clean

COUNT=bin/countwords.py

# Regenerate results for "Moby Dick"
results/moby_dick.csv : data/moby_dick.txt ${COUNT}
    python ${COUNT} data/moby_dick.txt > results/moby_dick.csv

# Regenerate results for "Jane Eyre"
results/jane_eyre.csv : data/jane_eyre.txt ${COUNT}
    python ${COUNT} data/jane_eyre.txt > results/jane_eyre.csv

# ...clean target as before...
```

The definition takes the form `NAME=value`. By convention, variables are written in UPPER CASE so that they'll stand out from filenames (which are usually in lower case), but that's not required. What *is* required is using `${NAME}` to refer to the variable: if we write `$NAME`, Make interprets that as “the variable called N followed by the three literal characters ‘AME’.” If no variable called N exists, `$NAME` becomes `AME`, which is almost certainly not useful.

Using variables doesn't just cut down on typing. They also make rules easier to understand, since they signal to readers that several things are always and exactly the same.

## 14.10 How can I make one update depend on another?

We can re-create all the results files with a single command by listing multiple targets when we run Make:

```
$ make results/moby_dick.csv results/jane_eyre.csv
```

However, users have to know what files they might want to create in order to do this, and have to type their names exactly right. A better approach is to create a phony target that depends on all of the output files and make it the first rule in the file so that it is the default. By convention, this target is called `all`, and while we don't have to list all our phony targets in alphabetical order, it makes them a lot easier to find:

```
.PHONY: all clean

COUNT=bin/countwords.py

# Regenerate all results.
all : results/moby_dick.csv results/jane_eyre.csv
```

```
# ...rules for moby_dick, jane_eyre, and clean...
```

If we run Make now, it sees that `all` is only “up to date” if the two CSV files are up to date, so it looks for a rule for each and runs each of those rules.

We can draw the prerequisites defined in the Makefile as a dependency graph, with arrows showing what each target depends on.

FIXME: figure

Note that the Makefile doesn’t define the order in which `results/moby_dick.csv` and `results/jane_eyre.csv` are updated, so Make can rebuild them in whatever order it wants. This is called declarative programming: we declare what outcome we want, and the program figures out how to achieve it.

## 14.11 How can I abbreviate my update rules?

We could add a third book to our Makefile, then a fourth, but any time we find ourselves duplicating code, there’s almost certainly a way to write a general rule. In order to create these, though, we first need to learn about automatic variables.

The first step is to use the very cryptic expression `$@` in the rule’s action to mean “the target of the rule”. We start with this:

```
# Regenerate results for "Moby Dick"
results/moby_dick.csv : data/moby_dick.txt
    python bin/countwords.py data/moby_dick.txt > results/moby_dick.csv
```

and turn it into this:

```
# Regenerate results for "Moby Dick"
results/moby_dick.csv : data/moby_dick.txt
    python bin/countwords.py data/moby_dick.txt > $@
```

`$@` is an automatic variable: Make defines its value for us separately in each rule. And yes, `$@` is an unfortunate name: something like `$TARGET` would be easier to understand, but we’re stuck with it.

Step 2 is to replace the list of prerequisites in the action with `$^`, which is another automatic variable meaning “all the prerequisites of the current rule”:

```
# Regenerate results for "Jane Eyre"
results/jane_eyre.csv : data/jane_eyre.txt
    python bin/countwords.py $^ > $@
```

But wait: our results files don’t just have books as dependencies. They also depend on `bin/countwords.py`. What happens if we include that in the rule

while using automatic variables? (We'll do this for a third book to keep the three rules separate in the example Makefile.)

```
# Regenerate results for "The Time Machine" - WRONG
results/time_machine.csv : data/time_machine.txt ${COUNT}
    python bin/countwords.py $^ > $@
```

This doesn't do the right thing because `$^` includes *all* of the prerequisites, so the action tries to process the script as if it were a data file:

```
python bin/countwords.py data/time_machine.txt bin/countwords.py results/time_machine.csv
```

This situation comes up so often that Make helpfully provides another automatic variable `$<` meaning “the first prerequisite”, which lets us rewrite our rules like this:

```
# Regenerate results for "Janey Eyre"
results/jane_eyre.csv : data/jane_eyre.txt ${COUNT}
    python bin/countwords.py $< > $@
```

And yes, `$< > $@` is hard to read, even with practice, and `< $<` (reading the first prerequisite from standard input) is even harder. Using an editor that does syntax highlighting helps (Chapter ??), and if you are ever designing software for other people to use, remember this case and don't do it.

## 14.12 How can I write one general rule to update many files in the same way?

We can now replace all the rules for generating results files with one pattern rule that uses `%` as a wildcard. Whatever part of a filename `%` matches in the target, it must also match in the prerequisites, so the single rule:

```
results/%.csv : data/%.txt ${COUNT}
    python bin/countwords.py $< > $@
```

will handle *Jane Eyre*, *Moby Dick*, and *The Time Machine*. (Unfortunately, `%` cannot be used in rules' actions, which is why `$<` and `$@` are needed.) With this rule in place, our entire Makefile is reduced to:

```
.PHONY: all clean
```

```
COUNT=bin/countwords.py
```

```
# Regenerate all results.
```

```
all : results/moby_dick.csv results/jane_eyre.csv results/time_machine.csv
```

```
# Regenerate result for any book.
```

```
results/%.csv : data/%.txt ${COUNT}
```

```
python ${COUNT} $< > $@

# Remove all generated files.
clean :
    rm -f results/*.csv

Let's delete all of the results files and re-create them all:

$ make -f pattern_rule.mk clean

rm -f results/*.csv

$ make -f pattern_rule.mk all

python bin/countwords.py data/moby_dick.txt > results/moby_dick.csv
python bin/countwords.py data/jane_eyre.txt > results/jane_eyre.csv
python bin/countwords.py data/time_machine.txt > results/time_machine.csv

We can still rebuild individual files:

$ touch data/jane_eyre.txt
$ make -f pattern_rule.mk results/jane_eyre.csv

python bin/countwords.py data/jane_eyre.txt > results/jane_eyre.csv
```

## 14.13 How can I define sets of files automatically?

Our “automated” analysis is still not fully automated: If we add another book to **raw**, we have to remember to add it to the **all** target in the Makefile as well. Once again, we will fix this in steps.

To start, imagine that all the results files already exist, and we just want to update them. We can define a variable called **RESULTS** to be a list of all the results files using the same notation we’d use in the shell to match all the CSV files in the **results/** directory:

```
RESULTS=results/*.csv
```

and then make **all** depend on that:

```
# Regenerate all results.
all : ${RESULTS}
```

This works, but only for re-creating files: if a results file doesn’t already exist when we run Make, its name won’t be included in **RESULTS**, and Make won’t realize that we want to generate it.

What we really want to do is generate the list of results files from the list of books in the **data/** directory. We can use a function to do this. The syntax is a

little odd, because functions were added to Make long after it was first written, but at least they have readable names. Let’s create a variable `DATA` that holds the names of all of our data files:

```
DATA = $(wildcard data/*.txt)
```

This calls the function `wildcard` with the argument `data/*.txt`. The result is a list of all the text files in the `raw` directory, just as we’d get with `data/*.txt` in the shell. (We could use a shell wildcard here as we did when defining `RESULTS`, but we want to show how functions work.)

Did this do the right thing? To check, we can add another phony target to the end of the file called `settings` that uses the shell command `echo` to print the name and value of a variable:

```
.PHONY: all clean settings

# ...everything else...

# Show variables' values.
settings :
    echo COUNT: ${COUNT}
    echo DATA: ${DATA}
```

Let’s run this:

```
$ make -f function_wildcard.mk settings

echo COUNT: bin/countwords.py
COUNT: bin/countwords.py
echo DATA: data/common_sense.txt data/jane_eyre.txt data/life_of_frederick_douglass.txt
DATA: data/common_sense.txt data/jane_eyre.txt data/life_of_frederick_douglass.txt data/...
```

The output appears twice because Make shows us the command it’s going to run before running it. If we put `@` before the command, Make doesn’t show it before running it:

```
settings :
    @echo COUNT: ${COUNT}
    @echo DATA: ${DATA}

$ make -f function_wildcard.mk settings

COUNT: bin/countwords.py
DATA: data/common_sense.txt data/jane_eyre.txt data/life_of_frederick_douglass.txt data/...
```

We now have the names of our input files, but what we want is the names of the corresponding output files. Make has another function called `patsubst` (short for “**p**attern **s**ubstitution”) that uses the same kind of patterns used in rules to do exactly this:

```
RESULTS=$(patsubst data/*.txt,results/%.csv,${DATA})
```



`$(patsubst ...)` calls the pattern substitution function. The first argument is what to look for: in this case, a text file in the `raw` directory. As in a pattern rule, we use `%` to match the stem of the file's name, which is the part we want to keep.

The second argument is the replacement we want. Ours uses the stem matched by `%` to construct the name of a CSV file in the `results` directory. Finally, the third argument is what we're doing substitutions in, which is our list of books' names.

Let's check that this has worked by adding to the `settings` target

```
settings :
    @echo COUNT: ${COUNT}
    @echo DATA: ${DATA}
    @echo RESULTS: ${RESULTS}

$ make -f patsubst.mk settings

COUNT: bin/countwords.py
DATA: data/common_sense.txt data/jane_eyre.txt data/life_of_frederick_douglass.txt data/moby_dick.txt
RESULTS: results/common_sense.csv results/jane_eyre.csv results/life_of_frederick_douglass.csv results/moby_dick.csv
```

Excellent: `DATA` has the names of all of the files we want to process and `RESULTS` automatically has the corresponding names of the files we want to generate. Let's test it:

```
$ make -f patsubst.mk clean

rm -f results/*.csv

$ make -f patsubst.mk all

python bin/countwords.py data/common_sense.txt > results/common_sense.csv
python bin/countwords.py data/jane_eyre.txt > results/jane_eyre.csv
python bin/countwords.py data/life_of_frederick_douglass.txt > results/life_of_frederick_douglass.csv
python bin/countwords.py data/moby_dick.txt > results/moby_dick.csv
python bin/countwords.py data/sense_and_sensibility.txt > results/sense_and_sensibility.csv
python bin/countwords.py data/time_machine.txt > results/time_machine.csv
```

Our workflow is now just two steps: add a data file and run Make. As we'll see in Chapter ??, we can even automate the second half in some cases, but this is still a big improvement over running things manually, particularly as we start to add more steps (such as merging data files and generating plots).

## 14.14 How can I document my workflow?

Every well-behaved program can tell people how to use it. If we run `make --help`, for example, we get a (long) list of things Make can do for us.

But how can we document the workflow that our Makefile now embodies? One common choice is to provide a special target like `settings` that prints a description of available targets:

```
.PHONY: all clean help settings

# ...other definitions...

# Show help.
help :
    @echo "all : regenerate all out-of-date results files."
    @echo "results/*.csv : regenerate a particular results file."
    @echo "clean : remove all generated files."
    @echo "settings : show the values of all variables."
    @echo "help : show this message."
```

This is easy to set up and does the job, but once again its redundancy should worry us: the same information appears in both the comments on rules and the help, which means that authors have to remember to update the help when adding or changing rules.

A better approach, which we will explore in more depth in Chapter ??, is to have people format some comments in a special way and then extract and display those comments when asked for help. We'll use `##` (a double comment marker) to indicate the lines we want displayed and use `grep` to extract lines that start with that marker. We will use Make's `MAKEFILE_LIST` variable to get the path to the Makefile, since we may be using the `-f` flag to specify which Makefile we're using. With all that in place, our finished Makefile is:

```
.PHONY: all clean help settings

COUNT=bin/countwords.py
DATA=$(wildcard data/*.txt)
RESULTS=$(patsubst data/%.txt,results/%.csv,{DATA})

## all : regenerate all results.
all : ${RESULTS}

## results/*.csv : regenerate result for any book.
results/%.csv : data/%.txt ${COUNT}
    python ${COUNT} $< > $@

## clean : remove all generated files.
clean :
    rm -f results/*.csv

## settings : show variables' values.
settings :
```

```

@echo COUNT: ${COUNT}
@echo DATA: ${DATA}
@echo RESULTS: ${RESULTS}

## help : show this message.
help :
    @grep '^##' ${MAKEFILE_LIST}

```

Let's test:

```
$ make -f makefile_grep.mk
```

```

## all : regenerate all results.
## results/*.csv : regenerate result for any book.
## clean : remove all generated files.
## settings : show variables' values.
## help : show this message.

```

With a bit more work we could strip off the leading `##` markers, but this is a good start.

### How did you know that?

FIXME: keep this personal?

I had never used the variable `MAKEFILE_LIST` before writing this lesson. In fact, until about 15 minutes ago, I didn't even know it existed: I always had my `help` target's action `grep` for `##` in `Makefile`. Once I realized that wouldn't work in this example (because I'm writing lots of little Makefiles to demonstrate ideas step by step) I searched online for "how to get name of Makefile in make". The second hit took me to the GNU Make documentation for other special variables, which told me exactly what I needed. I spend anywhere from a quarter to three quarters of my time searching for things when I program these days; one of the goals of these lessons is to give you an idea of what you ought to be searching for yourself so that you can do this more efficiently.

## 14.15 Summary

- Smith (2011) describes the design and implementation of several build tools in detail.

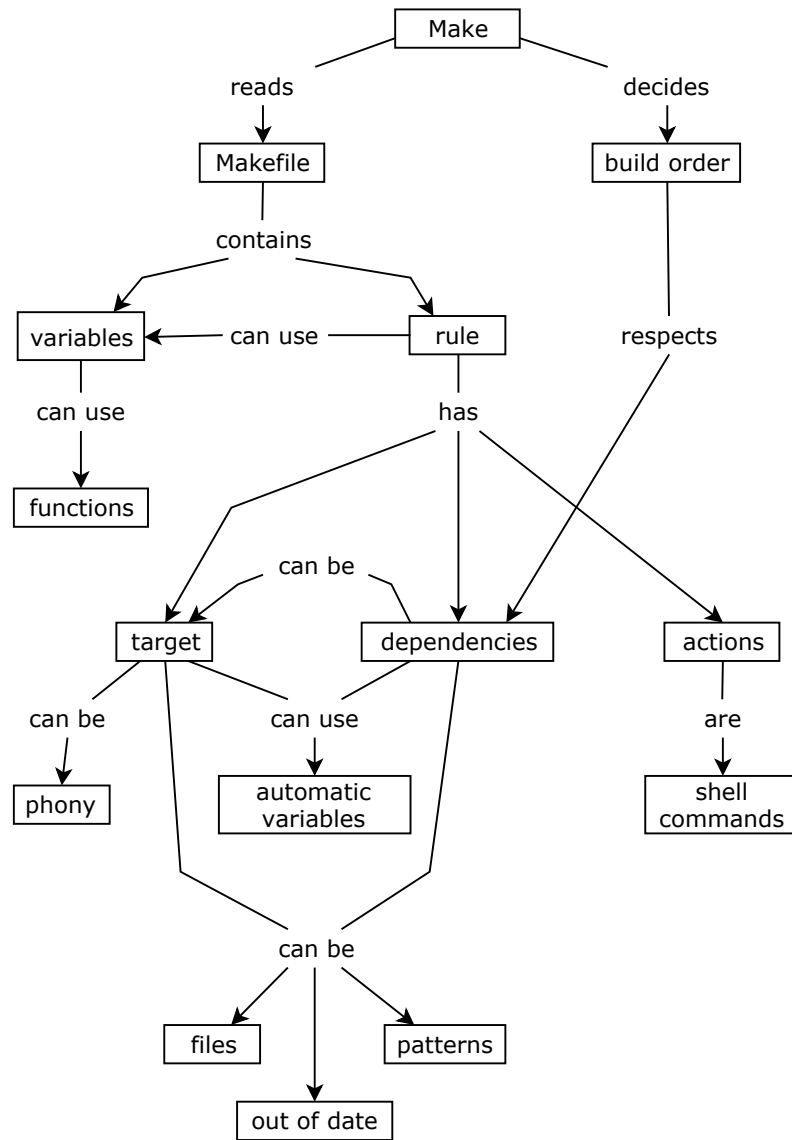


Figure 14.1: Automation Concept Map

## 14.16 Exercises

### 14.16.1 Create a summary CSV file

- Add rule to Makefile to create a summary CSV file from all of the book CSV files.
- Careful how to write the dependencies so that it doesn't depend on itself.

### 14.16.2 Generate a plot for the top N words

- FIXME: make it depend on the summary.

## 14.17 Key Points

- A build tool re-runs commands to bring files up to date with respect to the things they depend on.
- Make is a widely-used build tool that uses files' timestamps to find out-of-date prerequisites.
- A Make rule has targets, prerequisites, and actions.
- A target can correspond to a file or be a phony target (used simply to trigger actions).
- When a target is out of date with respect to its prerequisites, Make executes the actions associated with its rule.
- Make executes as many rules as it needs to when updating files, but always respect prerequisite order.
- Make defines the automatic variables `$@` (target), `$^` (all prerequisites), and `$<` (first prerequisite).
- Pattern rules can use `%` as a placeholder for parts of filenames.
- Makefiles can define variables using `NAME=value`.
- Makefiles can also use functions such as `$(wildcard ...)` and `$(patsubst ...)`.
- Specially-formatted comments can be used to make Makefiles self-documenting.



## Chapter 15

# Unit Testing

### 15.1 Questions

- How should I write tests for my software?
- How can I tell how much testing I've actually done?

### 15.2 Objectives

- Explain what realistic technical and social goals for software testing are.
- Explain what a test runner is.
- Explain what a text fixture is.
- Write and run unit tests using Python's `pytest` test runner.
- Check test coverage using Python's `coverage` module.

### 15.3 Introduction

Term frequency-inverse document frequency (TF-IDF) is a way to determine how relevant a document is in a search. A term's frequency is how often it occurs in a document  $d$  divided by the number of words in that document; inverse document frequency is the ratio of the total number of documents to the number of documents in which the word occurs:

$$\tau_w = \frac{n_w(d)/n_*(d)}{D/D_w}$$

If a word is very common in one document, but absent in others, TF-IDF is high for that (word, document) pair. If a word is common in all documents, or if a word is rare in a document, TF-IDF is lower. This simple measure can therefore be used to rank documents' relevance with respect to that term.

Calculating TF-IDF depends on counting words in a document... while ignoring leading/trailing punctuation... except for words like “Ph.D.” And then there's the question of hyphenation, which can either signal a multi-part word or a line break, unless of course it indicates that a multi-part happened to land at the end of a line. In short, this relatively simple description can be implemented in many subtly different ways, many of which are wrong. In order to tell which one we've actually built and whether we've built it correctly, we need to write some tests. This lesson therefore introduces some key ideas and tools.

This material is drawn in part from Testing and Continuous Integration with Python.

## 15.4 What are realistic goals for testing?

Any discussion of software and correctness has to start with two ideas. The first is that no amount of testing can ever prove that a piece of code is correct. A function that takes three arguments, each of which can have a hundred different values, theoretically needs a million tests. Even if we were to write them, how would we check that the values we're testing against are correct? And how would we tell that we had typed in those comparison values correctly, or that we were using the right equality checks in our code, or—the list goes on and on.

The second big idea is that the real situation is far less dire. Suppose we want to test a function that returns the sign of a number. If it works for the number 3, it will almost certainly work for 4, and for 5, and so on. In fact, there are only a handful of cases that we actually need to test:

| Value         | Expected | Reason                              |
|---------------|----------|-------------------------------------|
| <b>-Inf</b>   | -1       | Only value of its kind              |
| <b>Inf</b>    | 1        | Only value of its kind              |
| <b>NaN</b>    | NaN      | Only value of its kind              |
| <b>-5</b>     | -1       | Or some other negative number       |
| <b>0</b>      | 0        | The only value that has a sign of 0 |
| <b>127489</b> | 1        | Or some other positive number       |

As this table shows, we can divide test inputs into equivalence classes and check one member of each. Yes, there's a chance that we'll miss something—that the code we're testing will behave differently for values we have put in the same class—but this approach reduces the number of tests we have to write *and* makes



it easier for the next person reading our code to understand what it does.

## 15.5 What does a systematic software testing framework look like?

A framework for software testing has to:

- make it easy for people to write tests (because otherwise they won't do it);
- run a set of tests;
- report which ones have failed;
- give some idea of where or why they failed (to help debugging); and
- give some idea of whether any results have changed since the last run.

Any single test can have one of three results:

- success, meaning that the test passed correctly;
- failure, meaning that the software being tested didn't do what it was supposed to; or
- error, meaning that the test itself failed (in which case, we don't know anything about the software being tested).

A unit test is a function that runs some code and produces one of these three results. The input data to the unit test is called the fixture; we tell if the test passed or not by comparing the actual output to the expected output. For example, here's a very badly written version of `numSign` and an equally badly written pair of tests for it:

```
numSign <- function(x) {
  if (x > 0) {
    1
  } else {
    -1
  }
}

stopifnot(numSign(1) == 1)
stopifnot(numSign(-Inf) == -1)
stopifnot(numSign(NULL) == 0)
```

Here, the fixtures are 1, NULL, and `-Inf`, and the corresponding expected outputs are 1, 0, and -1. These tests are badly written for two reasons:

1. Execution halts at the first failed test, which means we get less information than we could about the state of the system we're testing.
2. Each test prints its output to the screen: there is no overall summary and no easy way to tell which test produced which result. This isn't a problem

when there are only three tests, but experience shows that if the output isn't comprehensible at a glance, developers will stop paying attention to it.

One other mistake that's often made in testing is making tests depend on each other. Good tests are independent: they should produce the same results no matter what other tests are run in what order, so that a failure in an early tests doesn't cause a later test to fail when it should pass or pass when it should fail. This means that every test should start from a freshly-generated fixture rather than using the output of previous tests.

## 15.6 How can I manage tests systematically?

A test framework is a library that provides us with functions that help us write tests, and includes a test runner that will find tests, execute them, and report both individual results that require attention and a summary of overall results.

There are many test frameworks for Python. One of the most popular is `pytest`, which structures tests according to three simple rules:

1. All tests are put in files whose names begin with `test_`.
2. Each test is a function whose name also begins with `test_`.
3. Test functions use `assert` to check that results are as expected.

For example, we could test the `count_words` function by putting the following code in `test_count.py`:

```
from tf_idf import count_word

def test_single_word():
    assert count_words('word') == {'word' : 1}

def test_single_repeated_word():
    assert count_words('word word') == {'word' : 2}

def test_two_words():
    assert count_words('another word') == {'another' : 1, 'word' : 1}

def test_trailing_punctuation():
    assert count_words("another's word") == {'another's' : 1, 'word' : 1}
```

The fixture in the last test is the string `"another's word"`, and the expected result is the dictionary `{'another's' : 1, 'word' : 1}`. Note that the `assert` statement doesn't include an error message; `pytest` will include the name of test function in its output if the test fails, and that name should be all the documentation we need. (If the test is so complicated that more is needed, we should probably write a simpler test.)

## 15.7. HOW CAN I TELL IF MY SOFTWARE FAILED AS IT WAS SUPPOSED TO?267

We can run our tests from the shell with a single command:

```
$ pytest
```

As it runs tests, pytest prints . for each one that passes and F for each one that fails. After all of the tests have been run, it prints a summary of the failures:

```
===== test session starts =====
platform darwin -- Python 3.6.5, pytest-3.5.1, py-1.5.3, pluggy-0.6.0
rootdir: /Users/pterry/still-magic/src/unit, inifile:
plugins: remotedata-0.2.1, openfiles-0.3.0, doctestplus-0.1.3, arraydiff-0.2
collected 4 items

test_count.py ...F [100%]

===== FAILURES =====
----- test_trailing_punctuation -----

    def test_trailing_punctuation():
>     assert count_words("anothers' word") == {'anothers' : 1, 'word' : 1}
E     assert {"anothers": 1, 'word': 1} == {'anothers': 1, 'word': 1}
E       Omitting 1 identical items, use -vv to show
E       Left contains more items:
E       {"anothers": 1}
E       Right contains more items:
E       {'anothers': 1}
E       Use -v to get the full diff

test_count.py:13: AssertionError
===== 1 failed, 3 passed in 0.05 seconds =====
```

pytest searches for all files named `test_*.py` or `*_test.py` in the current directory and its sub-directories. We can use command-line options to narrow the search: for example, `pytest test_count.py` runs only the tests in `test_count.py`. It automatically records and reports pass/fail/error counts and gives us a nicely-formatted report, but most importantly, it works the same way for everyone, so we can test without having to think about *how* (only about *what*). That said, fitting tests into this framework sometimes requires a few tricks, which we will explore in the sections that follow.

## 15.7 How can I tell if my software failed as it was supposed to?

Many errors in production systems happen because people don't test their error handling code. Yuan et al. (2014) found that almost all (92%) of catastrophic

system failures were the result of incorrect handling of non-fatal errors explicitly signalled in software, and that in 58% of the catastrophic failures, the underlying faults could easily have been detected through simple testing of error handling code. Our tests should therefore check that the software fails as it's supposed to and when it's supposed to; if it doesn't, we run the risk of a silent error.

We can check for exceptions manually using this pattern:

```
# Expect count_words to raise ValueError for empty input.
def test_text_not_empty():
    try:
        count_words('')
        assert False, 'Should not get this far'
    except ValueError:
        pass
```

This code runs the test (i.e., calls `count_words`) and then fails on purpose if execution reaches the next line. If the right kind of exception is raised, on the other hand, it does nothing. (We'll take a look in the exercises at how to improve this to catch the case where the wrong kind of exception is raised.)

This pattern works, but it violates our first rule of testing: if writing tests is clumsy, developers won't do it. To make life easier, `pytest` provides a context manager called `pytest.raises` to handle tests for exceptions. A context manager creates an object that lives exactly as long as a block of code, and which can do setup and cleanup actions at the start and end of that block. We can use `pytest.raises` with Python's `with` keyword to say that we expect a particular exception to be raised in that code block, and should report an error if one isn't raised:

```
import pytest

def test_text_not_empty():
    with pytest.raises(ValueError):
        count_words('')
```

```
===== test session starts =====
platform darwin -- Python 3.6.5, pytest-3.5.1, py-1.5.3, pluggy-0.6.0
rootdir: /Users/gvwilson/merely-useful/still-magic/src/unit, inifile:
plugins: remotedata-0.2.1, openfiles-0.3.0, doctestplus-0.1.3, arraydiff-0.2
collected 1 item
```

```
test_exception.py F [100%]
```

```
===== FAILURES =====
_____ test_text_not_empty _____

def test_text_not_empty():
```

```

        with pytest.raises(ValueError):
>         count_words('')
E         Failed: DID NOT RAISE <class 'ValueError'>

test_exception.py:6: Failed
===== 1 failed in 0.04 seconds =====

```

The output tells us that `count_words` doesn't raise an exception when given an empty string, so we should either decide that the count in this case is zero, or go back and fix our function.

## 15.8 How can I test software that includes randomness?

Data scientists use a lot of random numbers; testing code that relies on them makes use of the fact that they aren't actually random. A pseudorandom number generator (PRNG) uses a complex algorithm to create a stream of values that have all the properties of a truly random sequence. A key feature of PRNGs is that they can be initialized with a seed, and that given the same seed, the PRNG will always produce the same sequence of values. If you want your work to be reproducible, you should always seed your PRNG, and always record the seed somewhere so that you can re-run exactly the same calculations.

Pseudorandom numbers aren't the only unpredictable things in programs. If you rely on the current date and time, that counts as “randomness” as well. For example, suppose you want to test that a function correctly counts the number of dates between the start of the year and the current date. As time goes by, the correct answer will change, so how can you write a reusable function?

The answer is to write a wrapper for the function in question that either calls the actual function or does what you need for testing, and then to use that wrapper, and *only* that wrapper, everywhere in your program. The example below uses this approach: if the value `TESTING_DATE` has been set, `weeks_since_01` returns the number of weeks from `start` to that date, and if not, it returns the number of weeks from `start` to the current (unpredictable) date:

```

import datetime

DAYS_PER_WEEK = 7

TESTING_DATE = None
def weeks_since_01(start):
    current = TESTING_DATE
    if current is None:
        current = datetime.date.today()

```

```
return round((current - start).days / DAYS_PER_WEEK)
```

If this code is in a file called `wrappers.py`, we would use it like this:

```
print('current', wrappers.weeks_since_01(datetime.date(2018, 8, 1)))
wrappers.TESTING_DATE = datetime.date(2018, 8, 30)
print('fixed', wrappers.weeks_since_01(datetime.date(2018, 8, 1)))
```

```
current 25
fixed 4
```

A cleaner approach is to make the test control an attribute of the function. (Remember, functions are objects in memory like anything else, so we can attach other data to them.) Using this approach allows us to import the function on its own as a self-contained unit, and avoids making those functions depend on externally-defined (global) variables. The method works because the body of a function isn't executed as the function is defined, so it's OK to refer to values that are added afterward.

Here's what defining a wrapper function looks like:

```
def weeks_since_02(start):
    current = weeks_since_02.testing_date
    if current is None:
        current = datetime.date.today()
    return round((current - start).days / DAYS_PER_WEEK)
weeks_since_02.testing_date = None
```

and here's how we would use it:

```
# demo_test_weeks.py
print('first', weeks_since_02(datetime.date(2018, 8, 1)))
weeks_since_02.testing_date = datetime.date(2018, 8, 30)
print('second', weeks_since_02(datetime.date(2018, 8, 1)))
```

```
current 25
fixed 4
```

## 15.9 How can I test software that does I/O?

A lot of early books on unit testing said that tests shouldn't rely on external files, both because file I/O was slow and because those files could easily get lost. Neither stricture applies today: file I/O is much faster than it was in the 1990s, and if files of test data are stored in version control, they're no more likely to be lost than the source code.

That said, it's often easier to read unit tests if the "files" used as fixtures are

included right next to the tests, and it's easier to test *output* if the “files” that are created only ever live in memory: temporary output files can always be cleaned up after tests complete, but it's one extra burden on test authors.

A good way to avoid all of these problems is to treat strings in memory as if they were files, which is what Python's `StringIO` module does. A `StringIO` object has the same methods as a file object, but reads and writes text in memory instead of bytes on disk:

```
from io import StringIO

writer = StringIO()
for word in 'first second third'.split():
    writer.write('{}\n'.format(word))
print(writer.getvalue())
```

```
first
second
third
```

`StringIO` objects can also be initialized with some data that a program can read. (Notice that the lengths reported below include the newline character at the end of each line.)

```
DATA = '''first
second
third'''

for line in StringIO(DATA):
    print(len(line))
```

```
6
7
5
```

In order to use `StringIO` in tests, we may need to refactor our code a bit (Chapter ??). It's common to have a function open a file, read its contents, and return the result like this:

```
def main(infile, outfile):
    with open(infile, 'r') as reader:
        with open(outfile, 'w') as writer:
            # ...read from reader and write to writer...
```

However, this `main` function is hard to test, since there's no easy way to substitute a `StringIO` for the file inside that function. What we can do is reorganize the software so that file opening is done separately from reading and writing. This is good practice anyway for handling `stdin` and `stdout` in command-line tools, which don't need to be opened:

```
def main(infile, outfile):
    if infile == '-':
        reader = stdin
    else:
        reader = open(infile, 'r')
    if outfile == '-':
        writer = stdout
    else:
        writer = open(outfile, 'r')

    process(reader, writer)

    if infile == '-':
        reader.close()
    if outfile == '-':
        writer.close()

def process(reader, writer):
    # ...read from reader and write to writer...
```

After moving the reading and writing into `process`, we can easily pass in a couple of `StringIO` objects for testing.

## 15.10 How can I tell which parts of my software have (not) been tested?

Take a moment to study the code shown below. Can you tell which lines are and aren't being executed?

```
def first(left, right):
    if left < right:
        left, right = right, left
    while left > right:
        value = second(left, right)
        left, right = right, int(right/2)
    return value

def second(check, balance):
    if check > 0:
        return balance
    else:
        return 0

def main():
```



```

    final = first(3, 5)
    print(3, 5, final)

if __name__ == '__main__':
    main()

```

The answer is probably “no”, but the second half of the answer should be “that’s what computers are for”. Coverage measures which parts of program are and are not executed. In principle, a coverage tool keep a list of Booleans, one per line, all of which are initialized to **False**. Each time a line is executed, the coverage tool sets the corresponding flag to **True**. After the program finishes, the tool reports which lines have and have not been executed, along with summary statistics like the percentage of code executed.

It’s easy and wrong to obsess about meeting specific targets for test coverage. However, anything that *isn’t* tested should be assumed to be wrong, and drops in coverage often indicate new technical debt.

Use `pip install coverage` to install the standard Python coverage tool. Once you have done that, use `coverage run filename.py` instead of `python filename.py` to run your program. This creates a file in the current directory called `.coverage`; once your program completes, you can run `coverage report` to get a summary of the most recent report:

```

$ coverage run demo_coverage.py

Name                      Stmts   Miss  Cover
-----
demo_coverage.py          16      1    94%

```

If you want the details, you can use `coverage html` to generate an HTML listing:

FIXME: include unit/coverage.html

A more advanced tool called a profiler can give you even more information; we will discuss profilers briefly in Section ??.

## 15.11 Summary

One practice we haven’t described in this section is test-driven development (TDD). Rather than writing code and then writing tests, many developers believe we should write tests first to help us figure out what the code is supposed to do, and then write just enough code to make those tests pass. Once the code works, we should clean it up and commit it, then move on to the next task.

TDD’s advocates claim that writing tests first focuses people’s minds on what code is supposed to so that they’re not subject to confirmation bias when viewing

test results. They also claim that TDD ensures that code actually *is* testable, and that tests are actually written. However, the evidence backing these claims is contradictory: empirical studies have not found a strong effect Fucci et al. (2016), and at least one study suggests that it may not be the order of testing and coding that matters, but whether developers work in short, interleaved bursts Fucci et al. (2017). Many productive programmers still believe in TDD, so it's possible that we are measuring the wrong things.

FIXME: create concept map for unit testing

## 15.12 Exercises

### 15.12.1 Handle the wrong kind of exception

- FIXME: modify explicit exception testing code to handle the wrong kind of exception.

## 15.13 Key Points

- Testing can only ever show that software has flaws, not that it is correct.
- Its real purpose is to convince people (including yourself) that software is correct enough, and to make tolerances on 'enough' explicit.
- A test runner finds and runs tests written in a prescribed fashion and reports their results.
- A unit test can pass (work as expected), fail (meaning the software under test is flawed), or produce an error (meaning the test itself is flawed).
- A fixture is the data or other input that a test is run on.
- Every unit test should be independent of every other to keep results comprehensible and reliable.
- Programmers should check that their software fails when and as it is supposed to in order to avoid silent errors.
- Write test doubles to replace unpredictable inputs such as random numbers or the current date or time with a predictable value.
- Use string I/O doubles when testing file input and output.
- Use a coverage tool to check how well tests have exercised code.

# Chapter 16

## Verification

### 16.1 Questions

- How should I test a data analysis pipeline?

### 16.2 Objectives

- Explain why floating point results aren't random but can still be unpredictable.
- Explain why it is hard to test code that produces plots or other graphical output.
- Describe and implement heuristics for testing data analysis.
- Describe the role of inference in data analysis testing and use the `tda` library to find and check constraints on tabular data.

### 16.3 Introduction

The previous lesson explained how to test software in general; this one focuses on testing data analysis.

## 16.4 What is the difference between testing in software engineering and in data analysis?

Testing data analysis pipelines is often harder than testing mainstream software applications. The reason is that data analysts often don't know what the right answer is, which makes it hard to check correctness. The key distinction is the difference between validation and verification. Validation asks, "Is specification correct?" while verification asks, "It's the difference between building the right thing and building something right; the former question is often much harder for data scientists to answer."

Instead of unit testing, a better analogy is often physical experiments. When high school students are measuring acceleration due to gravity, they should get a figure close to

$$10m/sec^2$$

. Undergraduates might get

$$9.8m/sec^2$$

depending on the equipment used, but if either group gets

$$9.806m/sec^2$$

with a stopwatch, a marble, and an ramp, they're either incredibly lucky or cheating. Similarly, when testing data analysis pipelines, we often have to specify tolerances. Does the answer have to be exactly the same as a hand-calculated value or a previously-saved value? If not, how close is good enough?

We also need to distinguish between development and production. During development, our main concern is whether our answers are (close enough to) what we expect. We do this by analyzing small datasets and convincing ourselves that we're getting the right answer in some ad hoc way.

In production, on the other hand, our goal is to detect cases where behavior deviates significantly from what we previously decided what right. We want this to be automated so that our pipeline will ring an alarm bell to tell us something is wrong even if we're busy working on something else. We also have to decide on tolerances once again, since the real data will never have exactly the same characteristics as the data we used during development. We also need these checks because the pipeline's environment can change: for example, someone could upgrade a library that one of our libraries depends on, which could lead to us getting slightly different answers than we expected.

## 16.5 Why should I be cautious when using floating-point numbers?

Every tutorial on testing numerical software has to include a discussion of the perils of floating point, so we might as well get ours out of the way. The explanation that follows is simplified to keep it manageable; if you want to know more, please take half an hour to read Goldberg (1991).

Finding a good representation for floating point numbers is hard: we cannot represent an infinite number of real values with a finite set of bit patterns, and unlike integers, no matter what values we *do* represent, there will be an infinite number of values between each of them that we can't. These days, floating point numbers are usually represented using sign, magnitude (or mantissa), and an exponent. In a 32-bit word, the IEEE 754 standard calls for 1 bit of sign, 23 bits for the mantissa, and 8 bits for the exponent. To illustrate the problems with floating point, we will use a much simpler 5-bit representation with 3 bits for the magnitude and 2 for the exponent. We won't worry about fractions or negative numbers, since our simple representation will show off the main problems.

The table below shows the possible values (in decimal) that we can represent with 5 bits. Real floating point representations don't have all the redundancy that you see in this table, but it illustrates the point. Using subscripts to show the bases of numbers,

$$110_2 \times 2^{11_2}$$

$$6 \times 2^3$$

or 48.

Exponent

Mantissa

00

01

10

11

000

0

0

0

0

001

1

2

4

8

010

2

4

8

16

011

3

6

12

24

100

4

8

16

32

101

5

10

20

40

110

6

12

24

48

111

7



```
## 2 0.98999999999999991118 0.00000000000000000000
## 3 0.9989999999999999112 0.00000000000000000000
## 4 0.999900000000000011013 0.00000000000000000000
## 5 0.999990000000000045510 0.00000000000000000000
## 6 0.999999000000000082267 0.0000000000000000111022
## 7 0.999999900000000052636 0.00000000000000000000
## 8 0.999999990000000060775 0.0000000000000000111022
## 9 0.999999999000000028282 0.00000000000000000000
```

As the output shows, the very first value contributing to our sum is already slightly off. Even with 23 bits for a mantissa, we cannot exactly represent 0.9 in base 2, any more than we can exactly represent  $1/3$  in base 10. Doubling the size of the mantissa would reduce the error, but we can't ever eliminate it.

The good news is,

$$9 \times 10^{-1}$$

and

$$1 - 0.1$$

are exactly the same: the value might not be precisely right, but at least they are consistent. But some later values differ, and sometimes accumulated error makes the result *more* accurate.

It's very important to note that *this has nothing to do with randomness*. The same calculation will produce exactly the same results no matter how many times it is run, because the process is completely deterministic, just hard to predict. If you see someone run the same code on the same data with the same parameters many times and average the results, you should ask if they know what they're doing. (That said, doing this *can* be defensible if there is parallelism, which can change evaluation order, or if you're changing platform, e.g., moving computation to a GPU.)

## 16.6 How can I express how close one number is to another?

The absolute spacing in the diagram above between the values we can represent is uneven. However, the relative spacing between each set of values stays the same: the first group is separated by 1, then the separation becomes 2, then 4, and then 8. This happens because we're multiplying the same fixed set of mantissas by ever-larger exponents, and it leads to some useful definitions. The absolute error in an approximation is the absolute value of the difference between the approximation and the actual value. The relative error is the ratio of the absolute error to the value we're approximating. For example, if we are off by 1 in approximating  $8+1$  and  $56+1$ , we have the same absolute error, but the relative error is larger in the first case than in the second. Relative error is



almost always more important than absolute error when we are testing software because it makes little sense to say that we're off by a hundredth when the value in question is a billionth.

## 16.7 How should I write tests that involved floating-point values?

Accuracy is how close your answer is to right, and precision is how close repeated measurements are to each other. You can be precise without being accurate (systematic bias), or accurate without being precise (near the right answer, but without many significant digits). Accuracy is usually more important than precision for human decision making, and a relative error of

$$10^{-3}$$

(three decimal places) is more than good enough for most data science because the decision a human being would make won't change if the number changes by 0.1%.

We now come to the crux of this lesson: if the function you're testing uses floating point numbers, what do you compare its result to? If we compared the sum of the first few numbers in `vals` to what it's supposed to be, the answer could be `False` even though we're doing nothing wrong. If we compared it to a previously calculated result that we had stored somehow, the match would be exact.

No one has a good generic answer to this problem because its root cause is that we're using approximations, and each approximation has to be judged in context. So what can you do to test your programs? If you are comparing to a saved result, and the result was saved at full precision, you could use exact equality, because there is no reason for the new number to differ. However, any change to your code, however small, could trigger a report of a difference. Experience shows that these spurious warnings quickly lead developers to stop paying attention to their tests.

A much better approach is to write a test that checks whether numbers are the same within some tolerance, which is best expressed as a relative error. In Python, you can do this with `pytest.approx`, which works on lists, sets, arrays, and other collections, and can be given either relative or absolute error bounds. To show how it works, here's an example with an unrealistically tight absolute bound:

```
from pytest import approx

for bound in (1e-15, 1e-16):
    vals = []
```

```

for i in range(1, 10):
    number = 9.0 * 10.0 ** -i
    vals.append(number)
    total = sum(vals)
    expected = 1.0 - (10.0 ** -i)
    if total != approx(expected, abs=bound):
        print('{:22.21f} {:2d} {:22.21f} {:22.21f}'.format(bound, i, total, expected))

```

```

9.999999999999999790978e-17  6 0.999999000000000082267 0.99999899999999971244
9.999999999999999790978e-17  8 0.999999900000000060775 0.99999989999999949752

```

This tells us that two tests pass with an absolute error of

$$10^{-15}$$

but fail when the bound is

$$10^{-16}$$

, both of which are unreasonably tight. (Again, think of physical experiments: an absolute error of

$$10^{-15}$$

is one part in a trillion, which only a handful of high-precision experiments have ever achieved.)

## 16.8 How can I test plots and other graphical results?

Testing visualizations is hard: any change to the dimension of the plot, however small, can change many pixels in a raster image, and cosmetic changes such as moving the legend up a couple of pixels will similarly generate false positives.

The simplest solution is therefore *not* to test the generated image, but to test the data used to produce it. Unless you suspect that the plotting library contains bugs, feeding it the correct data should produce the correct plot.

If you *do* need to test the generated image, the only practical approach is to compare it to a saved image that you have visually verified. `pytest-mpl` does this by calculating the root mean square (RMS) difference between images, which must be below a threshold for the comparison to pass. It also allows you to turn off comparison of text, because font differences can throw up spurious failures. As with choosing tolerances for floating-point tests, your rule for picking thresholds should be, “If images are close enough that a human being would make the same decision about meaning, the test should pass”

FIXME: example

Another approach is to save the plot in a vector format like SVG that stores the coordinates of lines and other elements as text in a structure similar to that of HTML. You can then check that the right elements are there with the right properties, although this is less rewarding than you might think: again, small changes to the library or to plotting parameters can make all of the tests fail by moving elements by a pixel or two. Vector-based tests therefore still need to have thresholds on floating-point values.

## 16.9 How can I test the steps in a data analysis pipeline during development?

We can't tell you how to test your math, since we don't know what math you're using, but we *can* tell you where to get data to test it with. The first method is subsampling: choose random subsets of your data, analyze it, and see how close the output is to what you get with the full dataset. If output doesn't converge as sample size grows, something is probably unstable—which is not necessarily the same as wrong. Instability is often a problem with the algorithm, rather than with the implementation.

If you do this, it's important that you select a random sample from your data rather than (for example) the first  $N$  records or every  $N$ 'th record. If there is any ordering or grouping in your data, those techniques can produce samples that are biased, which may in turn invalidate some of your tests.

FIXME: add an exercise that subsamples the Zipf data.

The other approach is to test with synthetic data. With just a few lines of code, you can generate uniform data (i.e., data having the same values for all observations), strongly bimodal data (which is handy for testing clustering algorithms), or just sample a known distribution. If you do this, you should also try giving your pipeline data that *doesn't* fit your expected distribution and make sure that something, somewhere, complains. Doing this is the data science equivalent of testing the fire alarm every once in a while.

For example, we can write a short program to generate data that conforms to Zips' Law and use it to test our analysis. Real data will be integers (since words only occur or not), and distributions will be fractional. We will use 5% relative error as our threshold, which we pick by experimentation: 1% excludes a valid correct value. The test function is called `is_zipf`:

```
from pytest import approx
```

```
RELATIVE_ERROR = 0.05
```

```
def is_zipf(hist):
    scaled = [h/hist[0] for h in hist]
    print('scaled', scaled)
    perfect = [1/(1 + i) for i in range(len(hist))]
    print('perfect', perfect)
    return scaled == approx(perfect, rel=RELATIVE_ERROR)
```

Here are three tests that use this function with names that suggest their purpose:

```
def test_fit_correct():
    actual = [round(100 / (1 + i)) for i in range(10)]
    print('actual', actual)
    assert is_zipf(actual)

def test_fit_first_too_small():
    actual = [round(100 / (1 + i)) for i in range(10)]
    actual[0] /= 2
    assert not is_zipf(actual)

def test_fit_last_too_large():
    actual = [round(100 / (1 + i)) for i in range(10)]
    actual[-1] = actual[1]
    assert not is_zipf(actual)
```

## 16.10 How can I check the steps in a data analysis pipeline in production?

An operational test is one that is kept in place during production to tell users if everything is still working as it should. A common pattern for such tests is to have every tool append information to a log (Chapter E.3.5) and then have another tool check that log file after the run is over. Logging and then checking makes it easy to compare values between pipeline stages, and ensures that there's a record of why a problem was reported. Some common operational tests include:

- Does this pipeline stage produce the same number of output records as input records?
- Or fewer if the stage is aggregating?
- If two or more tables are being joined, is the number of output records equal to the product of the number of input records?
- Is the standard deviation be smaller than the range of the data?
- Are there any NaNs or NULLs where there aren't supposed to be any?

To illustrate these ideas, here's a script that reads a document and prints one line per word:

```
import sys

num_lines = num_words = 0
for line in sys.stdin:
    num_lines += 1
    words = [strip_punctuation(w) for w in line.strip().split()]
    num_words += len(words)
    for w in words:
        print(w)
with open('logfile.csv', 'a') as logger:
    logger.write('text_to_words.py,num_lines,{}\n'.format(num_lines))
    logger.write('text_to_words.py,num_words,{}\n'.format(num_words))
```

Here's a complementary script that counts how often words appear in its input:

```
import sys

num_words = 0
count = {}
for word in sys.stdin:
    num_words += 1
    count[word] = count.get(word, 0) + 1
for word in count:
    print('{} {}'.format(word, count[word]))
with open('logfile.csv', 'a') as logger:
    logger.write('word_count.py,num_words,{}\n'.format(num_words))
    logger.write('word_count.py,num_distinct,{}\n'.format(len(count)))
```

Both of these scripts write records to `logfile.csv`. When we look at that file after a typical run, we see records like this:

```
text_to_words.py,num_lines,431
text_to_words.py,num_words,2554
word_count.py,num_words,2554
word_count.py,num_distinct,1167
```

We can then write a small program to check that everything went as planned:

```
# read CSV file into the variable data
check(data['text_to_words.py']['num_lines'] <= data['word_count.py']['num_words'])
check(data['text_to_words.py']['num_words'] == data['word_count.py']['num_words'])
check(data['word_count.py']['num_words'] >= data['word_count.py']['num_distinct'])
```

## 16.11 How can I infer and check properties of my data?

Writing tests for the properties of data can be tedious, but some of the work can be automated. In particular, the TDDA library can infer test rules from data, such as `age <= 100`, `Date` should be sorted ascending, or `StartDate <= EndDate`. The library comes with a command-line tool called `tdda`, so that the command:

```
$ tdda discover elements92.csv elements.tdda
```

infers rules from data, while the command:

```
tdda verify elements92.csv elements.tdda
```

verifies data against those rules. The inferred rules are stored as JSON, which is (sort of) readable with a bit of practice. Reading the generated rules is a good way to get to know your data, and modifying values (e.g., changing the maximum allowed value for `Grade` from the observed 94.5 to the actual 100.0) is an easy way to make constraints explicit:

```
"fields": {
  "Name": {
    "type": "string",
    "min_length": 3,
    "max_length": 12,
    "max_nulls": 0,
    "no_duplicates": true
  },
  "Symbol": {
    "type": "string",
    "min_length": 1,
    "max_length": 2,
    "max_nulls": 0,
    "no_duplicates": true
  },
  "ChemicalSeries": {
    "type": "string",
    "min_length": 7,
    "max_length": 20,
    "max_nulls": 0,
    "allowed_values": [
      "Actinoid",
      "Alkali metal",
      "Alkaline earth metal",
      "Halogen",
      "Lanthanoid",
```

```

        "Metalloid",
        "Noble gas",
        "Nonmetal",
        "Poor metal",
        "Transition metal"
    ]
},
"AtomicWeight": {
    "type": "real",
    "min": 1.007947,
    "max": 238.028913,
    "sign": "positive",
    "max_nulls": 0
},
...
}

```

We can apply these inferred rules to all elements using the `-7` flag to get pure ASCII output and the `-f` flag to show only fields with failures:

```
$ tdda verify -7 -f elements118.csv elements92.tdda
```

FIELDS:

```
Name: 1 failure  4 passes  type OK  min_length OK  max_length X  max_nulls OK  no_duplicates OK
```

```
Symbol: 1 failure  4 passes  type OK  min_length OK  max_length X  max_nulls OK  no_duplicates OK
```

```
AtomicWeight: 2 failures  3 passes  type OK  min OK  max X  sign OK  max_nulls X
```

```
...others...
```

SUMMARY:

```
Constraints passing: 57
```

```
Constraints failing: 15
```

Another way to use TDDA is to generate constraints for two datasets and then look at differences in order to see how similar the datasets are to each other. This is especially useful if the constraint file is put under version control.

## 16.12 Summary

FIXME: create concept map for verification

## 16.13 Exercises

FIXME: create exercises for verification

## 16.14 Key Points

- Programmers should use tolerances when comparing floating-point numbers (not just in tests).
- Test the data structures used in plotting rather than the plots themselves.
- Check that parametric or non-parametric statistics of data do not differ from saved values by more than a specified tolerance.
- Infer constraints on data and then check that subsequent data sets obey these constraints.



## Chapter 17

# Project Structure

### 17.1 Questions

- How should I organize the files and directories in my project?

### 17.2 Objectives

- Describe and justify Noble’s Rules for organizing projects.
- Explain the purpose of README, LICENSE, CONDUCT, and CITATION files.

### 17.3 Introduction

Project organization is like a diet: there is no such thing as “no diet”, just a good one or a bad one. Similarly, there is no such thing as “no project organization”: your project is either organized well or poorly. As with coding style (Chapter ??), small pieces in predictable places with readable names are easier to find and use than large chunks that vary from project to project and have names like `stuff`.

### 17.4 What are Noble’s Rules?

Noble (2009) described a way to organize small bioinformatics projects that is equally useful for other kinds of research computing. Each project is put in a separate Git repository, and the directories in the root of this repository are

organized according to purpose. The original specification included five top-level directories:

- The `./src/` directory (short for “source”) holds source code for programs written in languages like C or C++ that need to be compiled. Many projects don’t have this directory because all of their code is written in languages that don’t need compilation.
- Runnable programs go in `./bin/` (an old Unix abbreviation for “binary”, meaning “not text”). This includes the compiled and runnable versions of C and C++ programs, and also shell scripts, Python or R programs, and everything else that can be executed.
- Raw data goes in `./data/` and is never modified after being stored.
- Results are put in `./results/`. This includes cleaned-up data, figures, and everything else that can be rebuilt using what’s in `./bin/` and `./data/`. If intermediate results can be re-created quickly and easily, they might not be stored in version control, but anything that is included in a manuscript should be here.
- Finally, documentation and manuscripts go in `./doc/`.

Figure 17.1 below shows this layout for a project called **g-trans**. A few things to notice are:

- The documentation for the **regulate** script appears in the root of `./doc/`, while the paper for JCMB is stored in a sub-directory, since it contains several files.
- The `./src/` directory contains a Makefile to re-build the **regulate** program (Chapter 14). Some projects put the Makefile in the root directory, reasoning that since it affects both `./src/` and `./bin/`, it belongs above them both rather than in either one.
- There are several sub-directories underneath `./data/` and `./results/`, which we will discuss in the next section. Each of the sub-directories in `./results/` has its own Makefile, which will re-create the contents of that directory.

## 17.5 How should files and sub-directories be named?

While the directories in the top level of each project are organized by purpose, the directories within `./data/` and `./results/` are organized chronologically so that it’s easy to see when data was gathered and when results were generated. These directories all have names in ISO date format like `YYYY-MM-DD` to make it easy to sort them chronologically. This naming is particularly helpful when data and results are used in several reports.

At all levels, filenames should be easy to match with simple shell wildcards.

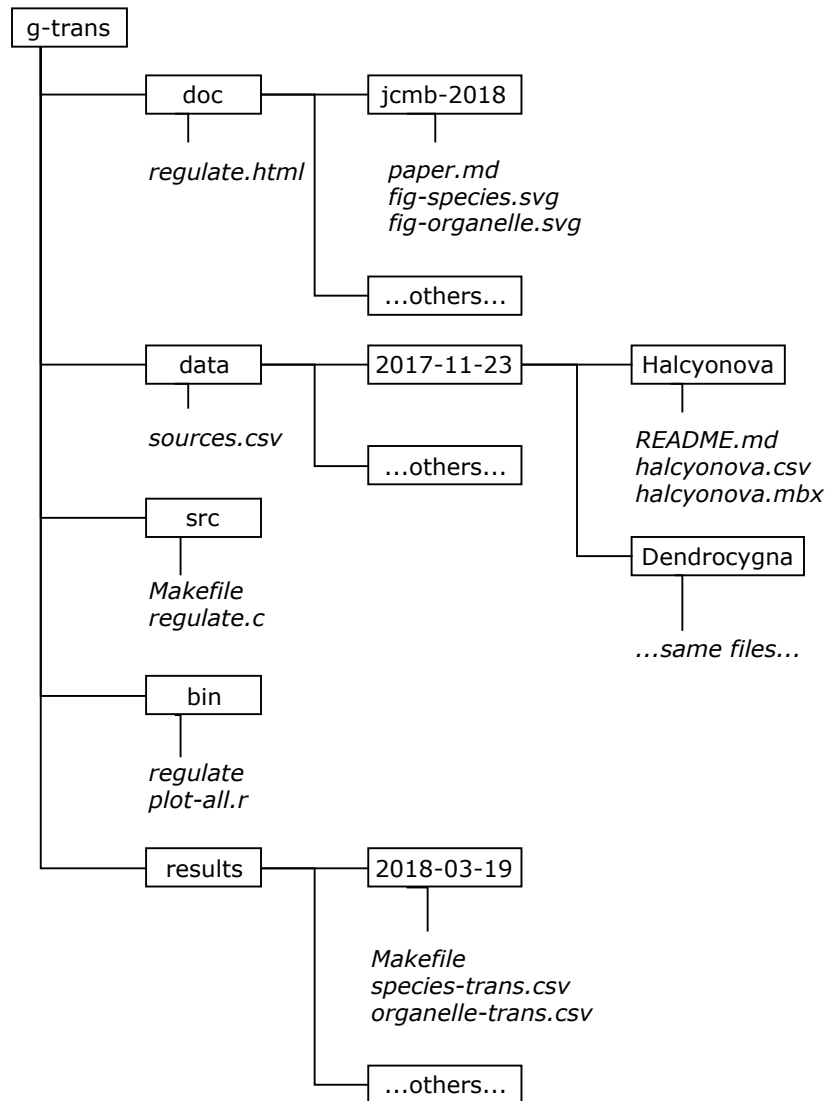


Figure 17.1: Project Layout

For example, a project might use `speciesorgantreatment.csv` as a file-naming convention, giving filenames like `human_kidney_cm200.csv`. This allows `human*_cm200.csv` to match all human organs or `*_kidney*.csv` to match all kidney data. It does produce long filenames, but tab completion means you only have to type them once. Long filenames are just as easy to match in programs: Python's `glob` and R's `Sys.glob` will both take a pattern and return a list of matching filenames.

## 17.6 How should I manage a mix of compiled programs and scripts?

Noble's Rules puts everything runnable in a single directory called `./bin/`. That makes things easy to find, but most software engineers would say that only the source code for programs should be in version control, not the output of the compiler. There are three ways to handle this:

1. Put the compiled program under version control. In theory this makes research more reproducible, since anyone who wants to re-run the analysis can be sure they're using exactly the same program as the author, but in practice, many "compiled" programs load libraries dynamically, so their results can still be affected by changes to their environment. Putting the compiled programs in version control *does* make it easier for people to re-run analyses, though, since they don't need to be able to compile things themselves.
2. Put everything in `./bin/`, then edit `.gitignore` to tell Git to ignore the compiled programs. This keeps everything runnable in one place, but requires a little bit of extra bookkeeping.
3. Put compiled programs in `./bin/` and everything that doesn't require compilation in `./scripts/`. This makes management simpler (just ignore everything in `./bin/`), but means there are two places where runnable programs might live.

We usually go with the third option, but they're all good as long as you are consistent between projects and document your rule.

## 17.7 Should I separate documentation from manuscripts?

Noble's Rules place documentation and manuscripts in `./docs/`. As with compiled programs and scripts, some people separate these, so that (for example) `./docs/` has the project's website and the documentation for its software,

while `./reports/` has one sub-directory for each paper, thesis chapter, or other manuscript. As more researchers [work in the open(#g:open-science), and as tools like R Markdown and the Jupyter Notebook blur the distinction between software, documentation, and reports, separating the two makes less sense.

## 17.8 How should I handle data can't be stored in version control?

Small datasets that don't contain sensitive information can and should be stored in version control. As a rule of thumb, anything that you would send in an email attachment is probably small enough to be put into Git, while anything that might reveal someone's identity when combined with other data sets should not be.

If data is large or sensitive, there should still be something in `./data/` to show its existence, and that "something" should be easy for programs to read. One option is a CSV file whose columns are:

- the name of the dataset,
- its URL or other unique identifier,
- the date it was last checked, and
- its size (so that users will have some idea of how much work is involved in processing it).

Another option is to have one file per dataset, so that instead of reading `human_genome.bam`, the program reads `human_genome.yml` and then uses the `url` key in that file to find the data it actually wants. If your data is complicated, you may also want to include a `README.md` file in `./data/` modelled on the one that accompanies The Pudding article "Women's Pockets are Inferior" (see Chapter 24).

## 17.9 What other files should every project contain?

Most projects' repositories contain a few files that weren't mentioned in Noble's paper, but which have become conventional in open source projects. All of these files may be plain text or Markdown, and may have a `.txt` or `.md` suffix (or no suffix at all), but please use the principal names given, in all caps, since a growing number of tools expect them.

- **README**: the project's title and a one-paragraph description of its purpose or content. This file is displayed by GitHub Pages as the project's home page.

- **LICENSE**: the project’s license (discussed in Chapter 18).
- **CONDUCT**: its code of conduct (also discussed in Chapter 18).
- **CITATION**: how the work should be cited. This file usually contains a plain text citation, and may also include entries formatted for various bibliographic systems like BibTeX. Some projects also have a separate **CONTRIBUTORS** file listing everyone who has contributed to it, while others include this as a section in **CITATION**.

If setting up the project or contributing to it are complex, there may also be a file called **BUILD** that explains how to install the required software, the formatting conventions for new data entries, and so on. These instructions can also be included as a section in **README**. Whichever convention is used, remember that the easier it is for people to get set up and contribute, the more likely they are to do so Steinmacher et al. (2014).

## 17.10 What *is* a project?

Like features (Chapter 12), what exactly constitutes a “project” requires a bit of judgment, and different people will make different decisions. Some common criteria are one project per publication, one project per deliverable piece of software, or one project per team. The first tends to be too small: a good data set will result in several reports, and the goal of some projects is to produce a steady stream of reports (such as monthly forecasts). The second is a good fit for software engineering projects whose primary aim is to produce tools rather than results, but is inappropriate for most data analysis work. The third tends to be too large: a team of half a dozen people may work on many different things at once, and a repository that holds them all quickly looks like someone’s basement.

The best rule of thumb we have found for deciding what is and isn’t a project is to ask what you have meetings about. If the same set of people need to get together on a regular basis to talk about something, that “something” probably deserves its own repository. And if the list of people changes slowly over time but the meetings continue, that’s an even stronger sign.

## 17.11 Summary

**FIXME**: create concept map for project structure

## 17.12 Exercises

FIXME: exercises for structure chapter.

## 17.13 Key Points

- Put source code for compilation in `./src/`.
- Put runnable code in `./bin/`.
- Put raw data in `./data/`.
- Put results in `./results/`.
- Put documentation and manuscripts in `./doc/`.
- Use file and directory names that are easy to match and include dates for the level under `./data/` and `./results/`.
- Create README, LICENSE, CONDUCT, and CITATION files in the root directory of the project.





## Chapter 18

# Including Everyone

### 18.1 Questions

- Why should I make my project welcoming for everyone?
- Why do I need a license for my work?
- What license should I use for my work?
- How should I tell people what they can and cannot do with my work?
- Why should my project have an explicit Code of Conduct?
- How can I be a good ally?

### 18.2 Objectives

- Explain the purpose of a Code of Conduct and the essential features an effective one must have.
- Explain why adding licensing information to a repository is important.
- Explain differences in licensing and social expectations.
- Choose an appropriate license.
- Explain where and how to communicate licensing.
- Explain steps a project lead can take to be a good ally.

### 18.3 Introduction

The previous lesson talked about the physical organization of projects. This one talks about the social structure, which is more important to the project's success. A project can survive badly-organized code; none will survive for long if people are confused, pulling in different directions, or hostile. This lesson

therefore talks about what projects can do to make newcomers feel welcome and to make things run smoothly after that. It draws on Fogel (2005), which describes how good open source software projects are run, and on Bollier (2014), which explains what a commons is and when it's the right model to use.

## 18.4 Why does a project need explicit rules?

Jo Freeman's influential essay "The Tyranny of Structurelessness" pointed out that every group has a power structure; the only question is whether it is formal and accountable or informal and unaccountable Freeman (1972). Thirty-five years after the free software movement took on its modern, self-aware form, its successes and failures have shown that if a project doesn't clearly state who has the right to do what, it will wind up being run by whoever argues loudest and longest.

## 18.5 How should I license my software?

It might seem strange to put licensing under discussion of inclusivity, but if the law or a publication agreement prevents people from reading your work or using your software, you're excluding them (and probably hurting your own career). You may need to do this in order to respect personal or commercial confidentiality, but the first and most important rule of inclusivity is to be open by default.

However, that is easier said than done, not least because the law hasn't kept up with everyday practice. Morin and Sliz (2012) and this blog post are good starting points from a scientist's point of view, while Lindberg (2008) is a deeper dive for those who want details. In brief, creative works are automatically eligible for intellectual property (and thus copyright) protection. This means that every creative work has some sort of license: the only question is whether authors and users know what it is.

Every project should therefore include an explicit license. This license should be chosen early: if you don't set it up right at the start, then each collaborator will hold copyright on their work and will need to be asked for approval when a license *is* chosen. By convention, the license is usually put in a file called `LICENSE` or `LICENSE.txt` in the project's root directory. This file should clearly state the license(s) under which the content is being made available; the plural is used because code, data, and text may be covered by different licenses.

*Don't write your own license*, even if you are a lawyer: legalese is a highly technical language, and words don't mean what you think they do.

To make license selection as easy as possible, GitHub allows you to select one of the most common licenses when creating a repository. The Open Source Initiative maintains a list of licenses, and [choosealicense.com](http://choosealicense.com) will help you find a license that suits your needs. Some of the things you will need to think about are:

1. Do you want to license the code at all?
2. Is the content you are licensing source code?
3. Do you require people distributing derivative works to also distribute their code?
4. Do you want to address patent rights?
5. Is your license compatible with the licenses of the software you depend on? For example, as we will discuss below, you can use MIT-licensed code in a GPL-licensed project but not vice versa.

The two most popular licenses for software are the MIT license and the GNU Public License (GPL). The MIT license (and its close sibling the BSD license) say that people can do whatever they want to with the software as long as they cite the original source, and that the authors accept no responsibility if things go wrong. The GPL gives people similar rights, but requires them to share their own work on the same terms:

You may copy, distribute and modify the software as long as you track changes/dates in source files. Any modifications to or software including (via compiler) GPL-licensed code must also be made available under the GPL along with build & install instructions.

— tl;dr

We recommend the MIT license: it places the fewest restrictions on future action, it can be made stricter later on, and the last thirty years shows that it's good enough to keep work open.

## 18.6 How should I license my data and reports?

The MIT license and the GPL apply to software. When it comes to data and reports, the most widely used family of licenses are those produced by Creative Commons, which have been written and checked by lawyers and are well understood by the community.

The most liberal license is referred to as CC-0, where the “0” stands for “zero restrictions”. CC-0 puts work in the public domain, i.e., allows anyone who wants to use it to do so however they want with no restrictions. This is usually the best choice for data, since it simplifies aggregate analysis. For example, if you choose a license for data that requires people to cite their source, then anyone who uses that data in an analysis must cite you; so must anyone who cites *their* results, and so on, which quickly becomes unwieldy.

The next most common license is the Creative Commons - Attribution license, usually referred to as CC-BY. This allows people to do whatever they want to with the work as long as they cite the original source. This is the best license to use for manuscripts, since you *want* people to share them widely but also want to get credit for your work.

Other Creative Commons licenses incorporate various restrictions on specific use cases:

- ND (no derivative works) prevents people from creating modified versions of your work. Unfortunately, this also inhibits translation and reformatting.
- NC (no commercial use) does *not* mean that people cannot charge money for something that includes your work, though some publishers still try to imply that in order to scare people away from open licensing. Instead, the NC clause means that people cannot charge for something that uses your work without your explicit permission, which you can give under whatever terms you want.
- Finally, SA (share-alike) requires people to share work that incorporates yours on the same terms that you used. Again, this is fine in principle, but in practice makes aggregation a headache.

## 18.7 Why should I establish a code of conduct for my project?

You don't expect to have a fire, but every large building or event should have a fire safety plan. Similarly, having a Code of Conduct for your project reduces the uncertainty that participants face about what is acceptable and unacceptable behavior. You might think this is obvious, but long experience shows that articulating it clearly and concisely reduces problems caused by have different expectations, particularly when people from very different cultural backgrounds are trying to collaborate. An explicit Code of Conduct is particularly helpful for newcomers, so having one can help your project grow and encourage people to give you feedback.

Having a Code of Conduct is particularly important for people from marginalized or under-represented groups, who have probably experienced harassment or unwelcoming behavior before. By adopting one, you signal that your project is trying to be a better place than YouTube, Twitter, and other online cesspools. Some people may push back claiming that it's unnecessary, or that it infringes freedom of speech, but in our experience, what they often mean is that thinking about how they might have benefited from past inequity makes them feel uncomfortable, or that they like to argue for the sake of arguing. If having a Code of Conduct leads to them going elsewhere, that will probably make your project run more smoothly.

Just as you shouldn't write your own license for a project, you probably shouldn't write your own Code of Conduct. We recommend using the Contributor Covenant for development projects and the model code of conduct from the Geek Feminism Wiki for in-person events. Both have been thought through carefully and revised in the light of experience, and both are now used widely enough that many potential participants in your project will not need to have them explained.

Rules are meaningless if they aren't enforced. If you adopt a Code of Conduct, it is therefore important to be clear about how to report issues and who will handle them. Aurora et al. (2018) is a short, practical guide to handling incidents; like the Contributor Covenant and the model code of conduct, it's better to start with something that other people have thought through and refined than to try to create something from scratch.

## 18.8 Why can I be a good ally for members of marginalized groups?

Setting out rules and handling incidents when they arise is what projects can do; if you have power (even or especially the power that comes from being a member of the majority group), what you can do personally is be a good ally for members of marginalized groups. Much of this discussion is drawn from the Frameshift Consulting Ally Skills workshop, which you should attend if you can.

First, some definitions. Privilege is an unearned advantage given to some people but not all; oppression is systemic, pervasive inequality that benefits the privileged and harms those without privilege. A straight, white, physically able, economically secure male is less likely to be interrupted when speaking, more likely to be called on in class, and more likely to get a job interview based on an identical CV than someone who is perceived as being outside these categories. The unearned advantage may be small in any individual case, but compound interest quickly amplifies these differences: someone who is called on more often in class is more likely to be remembered by a professor, who in turn is therefore more likely to recommend them to a potential employer, who is more likely to excuse the poor grades on their transcripts, and on and on it goes. People who are privileged are often not aware of it for the same reason that most fish don't know what water tastes like.

A target is someone who suffers from oppression. Targets are often called "members of a marginalized group", but that phrasing is deliberately passive. Targets don't choose to be marginalized: those with privilege marginalize them. An ally is a member of a privileged group who is working to understand their own privilege and end oppression. For example, privilege is being able to walk into a store and have the owner assume you're there to buy things, not to steal them. Oppression is the stories told about (for example) indigenous people being thieves,

and the actions people take as a result of them. A target is an indigenous person who wants to buy milk, and an ally is a white person who pays attention to a lesson like this one (raising their own awareness), calls out peers who spread racist stories (a peer action), or asks the shopkeeper whether they should leave too (a situational action).

Why should you be an ally? You could do it out of a sense of fairness because you realize that you have benefited from oppression even if you haven't taken part (or don't think you have). And you should do it because you can: taking action to value diversity results in worse performance ratings for minority and female leaders, while ethnic majority or male leaders who do this aren't penalized Hekman et al. (2017). As soon as you acknowledge that (for example) women are called on less often than men, or are less likely to get an interview or a publication given identical work, you have to acknowledge that white and Asian males are *more* likely to get these benefits than their performance alone deserves.

## 18.9 How can I be a good ally for members of marginalized groups?

So much for the theory: what should you actually do? A few simple rules will go a long way:

1. Be short, simple, and firm.
2. Don't try to be funny: it almost always backfires, or will later be used against you.
3. Play for the audience: you probably won't change the mind of the oppressor you're calling out, but you might change the minds or give heart to people who are observing.
4. Pick your battles. You can't challenge everyone, every time, without exhausting yourself and deafening your audience. An occasional sharp retort will be much more effective than constant criticism.
5. Don't shame or insult one group when trying to help another. For example, don't call someone stupid when what you really mean is that they're racist or homophobic.
6. Change the terms of the debate.

The last rule is best explained by example. Suppose someone says, "Why should we take diversity into account when hiring? Why don't we just hire the best candidate?" Your response could be, "Because taking diversity into account *is* hiring the best candidate. If you can run a mile in four minutes and someone else can do it in 4:15 with a ball and chain on their leg, who the better athlete? Who will perform better *if the impediment is removed*? If you intend to preserve an exclusionary culture in this lab, considering how much someone has achieved despite systemic unfairness might not make sense, but you're not arguing for that, are you?" And if someone then says, "But it's not fair to take anything

other than technical skill into account when hiring for a technical job,” you can say, “You’re right, which means that what you’re *really* upset about is the thought that you might be treated the way targets have been treated their whole lives.”

Captain Awkward has useful advice, and Charles’ Rules of Argument are very useful online:

1. Don’t go looking for an argument.
2. State your position once, speaking to the audience.
3. Wait for absurd replies.
4. Reply once more to correct any misunderstandings of your original statement.
5. Do not reply again—go do something fun instead.

Finally, it’s important to recognize that good principles sometimes conflict. For example, consider this scenario:

A manager consistently uses male pronouns to refer to software and people of unknown gender. When you tell them it makes you uncomfortable to treat maleness as the norm, they say that male is the default gender in their first language and you should be more considerate of people from other cultures.

On the one hand, you want to respect other people’s cultures; on the other hand, you want to be inclusive of women. In this case, the manager’s discomfort about changing pronouns matters less than the career harm caused by them being exclusionary, but many cases are not this clear cut. Like any written rules, a Code of Conduct requires constant interpretation; like everything else, discussion about specific cases becomes easier with practice.

## 18.10 Summary

FIXME: create concept map for making an inclusive project

## 18.11 Exercises

FIXME: exercises for creating an inclusive project.

## 18.12 Key Points

- Create an explicit Code of Conduct for your project modelled on the Contributor Covenant.

- Be clear about how to report violations of the Code of Conduct and who will handle such reports.
- People who are not lawyers should not try to write licenses.
- Every project should include an explicit license to make clear who can do what with the material.
- People who incorporate GPL'd software into their own software must make their software also open under the GPL license; most other open licenses do not require this.
- The Creative Commons family of licenses allow people to mix and match requirements and restrictions on attribution, creation of derivative works, further sharing, and commercialization.
- Be proactive about welcoming and nurturing community members.



## Chapter 19

# Managing Backlog

### 19.1 Questions

- How can I tell what needs to be done and who is doing it?

### 19.2 Objectives

- Explain what an issue tracking tool does and what it should be used for.
- Explain how to use labels on issues to manage work.
- Describe the information a well-written issue should contain.

### 19.3 Introduction

Version control tells us where we've been; issues tells us where we're going. Issue tracking tools are often called ticketing systems or bug trackers because they were created to keep track of work that needs to be done and bugs that needed fixing. However, they can be used to manage any kind of work and are often a convenient way to manage discussions as well.

### 19.4 How can I manage the work I still have to do?

Like other forges, GitHub records a set of issues for each project, and allows members of that project to create new ones, modify existing ones, and search

both. Every issue has:

- A unique ID, such as #123, which is also part of its link. This makes issues easy to find and refer to: in particular, GitHub automatically translates the expression #123 in a commit message into a link to that issue.
- A one-line title to aid browsing and search.
- The issue's current status. In simple systems (like GitHub's), each issue is either open or closed, and by default, only open issues are displayed.
- The ID of the issue's creator. The IDs of people who have commented on it or modified it are also embedded in the issue's history, which helps when using issues to manage discussions.
- A full description that may include screenshots, error messages, and just about anything else.
- Replies, counter-replies, and so on from people with an interest in the issue.

Broadly speaking, people create three kinds of issues:

1. *Bug reports* to describe problems they have encountered. A good bug report is a work of art, so we discuss them in detail below.
2. *Feature requests* (for software packages) or *tasks* (for projects). These issues are the other kind of work backlog in a project: not “what needs to be fixed” but “what do we want to do next”.
3. *Questions*. Many projects encourage people to ask questions on a mailing list or in a chat channel, but answers given there can be hard to find later, which leads to the same questions coming up over and over again. If people can be persuaded to ask questions by filing issues, and to respond to issues of this kind, then the project's old issues become a customized Stack Overflow for the project. (In practice, most projects find that it's easier to create a page of links to old questions and answers, or to rely on search to find things, since using issues this way requires someone to put in curation time.)

## 19.5 How can I write a good bug report?

The better the bug report, the faster the response, and the more likely the response will actually address the issue Bettenburg et al. (2008).

1. Make sure the problem actually *is* a bug. It's always possible that you have called a function the wrong way or done an analysis using the wrong configuration file; if you take a minute to double-check, you could well fix the problem yourself.

2. Try to come up with a reproducible example, or reprex. A reprex includes only the steps or lines of code needed to make the problem happen; again, you'll be surprised how often you can solve the problem yourself as you trim down your steps to create one.
3. Write a one-line title for the issue and a longer (but still brief) description that includes relevant details).
4. Attach any screenshots that show the problem or (slimmed-down) input files needed to re-create it.
5. Describe the version of the software you were using, the operating system you were running on, and anything else that might affect behavior.
6. Describe each problem separately so that each one can be tackled on its own. (This parallels the rule about creating a branch in version control for each bug fix or feature discussed in Chapter 12.)

Here's an example of a well-written bug report with all of the fields mentioned above:

```
ID: 1278
Creator: standage
Owner: malvika
Labels: Bug, Assigned
Summary: wordbase.py fails on accented characters
Description:
```

1. Create a text file called 'accent.txt' containing the word "Pumpernickel" with an umlaut over the 'u'.
2. Run 'python wordbase.py --all --message accent.txt'

Program should print "Pumpernickel" on a line by itself with the umlaut, but instead it fails with the message:

```
No encoding for [] on line 1 of 'accent.txt'.
```

([] shows where a solid black box appears in the output instead of a printable character.)

Versions:

- ``python wordbase.py --version`` reports 0.13.1
- Using on Windows 10.

## 19.6 How can I use labels to organize work?

There is always more work to do than there is time to do it. Issue trackers let project members add labels to issues to manage this. A label is just a word or

two; GitHub allows project owners to create any labels they want that users can then add to their issues. A small project should always use these four labels:

- *Bug*: something should work but doesn't.
- *Enhancement*: something that someone wants added.
- *Task*: something needs to be done, but won't show up in code (e.g., organizing the next team meeting).
- *Discussion*: something the team needs to make a decision about. (All issues can have discussion—this category is for issues that start that way.)

It might also use:

- *Question*: how is something supposed to work (discussed earlier).
- *Suitable for Newcomer* or *Beginner-Friendly*: to identify an easy starting point for someone who has just joined the project. If you help potential new contributors find places to start, they're more likely to do so Steinmacher et al. (2014).

## 19.7 How can I use labels to prioritize work?

The labels listed above described what kind of work an issue describes; a separate set of labels can be used to indicate how important that work is. These labels typically have names like “High Priority” or “Low Priority”; their purpose is to help triage, which is the process of deciding what is going to be worked on in what order.

In a large project, this is the responsibility of the product manager (Chapter ??); in a small one, it's common for the project's lead to decide this, or for project members to use up-votes and down-votes on an issue to indicate how important they think it is.

However the decision is made, it's helpful to use six more labels to record the outcomes:

- *Urgent*: this needs to be done *now* (typically reserved for security fixes).
- *Current*: this issue is included in the current round of work.
- *Next*: this issue is (probably) going to be included in the next round.
- *Eventually*: someone has looked at the issue and believes it needs to be tackled, but there's no immediate plan to do it.
- *Won't Fix*: someone has decided that the issue isn't going to be addressed, either because it's out of scope or because it's not actually a bug. Once an issue is marked this way, it is usually then closed.
- *Duplicate*: this issue is a duplicate of one that's already in the system. Again, issues marked this way are usually then closed.

Larger projects will replace “Current”, “Next”, and “Eventually” with labels corresponding to upcoming software releases, journal issues, or conferences. Un-

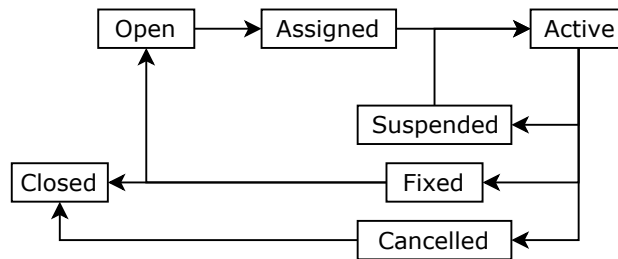


Figure 19.1: Issue State Transitions

fortunately, this doesn’t work well on GitHub’s issue tracking system, since there’s no way to “retire” a label without deleting it. After a while, a project that uses labels for specific events can have a *lot* of old labels...

## 19.8 How can I use issues to enforce a workflow for a project?

If you’re using GitHub, the short answer to this section’s title question is, “You can’t,” but you can create conventions, and more sophisticated issue tracking systems will let you enforce those conventions. The basic idea is to only allow issues to move from some states to others, and to only allow some people to move them. For example, Figure 19.1 shows a workflow for handling bugs in a medium-sized project:

- An “Open” issue becomes “Assigned” when someone is made responsible for it.
- An “Assigned” issue becomes “Active” when that person starts to work on it.
- If they stop work for any length of time, it becomes “Suspended”. (That way, people who are waiting for it know not to hold their breath.)
- An “Active” bug can either be “Fixed” or “Cancelled”, where the latter state means that the person working on it has decided it’s not really a bug.
- Once a bug is “Fixed”, it can either be “Closed” or, if the fix doesn’t work properly, moved back into the “Open” state and go around again.

Workflows like this are more complex than small projects need, but as soon as the team is distributed—particularly distributed across timezones—it’s very helpful to be able to find out what’s going on without having to wait for someone to be awake.

## 19.9 Summary

FIXME: create concept map for workflow

## 19.10 Exercises

FIXME: exercises for backlog chapter.

## 19.11 Key Points

- Create issues for bugs, enhancement requests, and discussions.
- Add people to issues to show who is responsible for working on what.
- Add labels to issues to identify their purpose.
- Use rules for issue state transitions to define a workflow for a project.

## Chapter 20

# Teamwork

### 20.1 Questions

- FIXME

### 20.2 Objectives

- FIXME

### 20.3 Introduction

FIXME: introduction

### 20.4 How can we run meetings more efficiently?

Most people do meetings poorly: they don't have an agenda going in, they don't take minutes, they waffle on or wander off into irrelevancies, they repeat what others have said or recite banalities simply so that they'll have said something, and they hold side conversations (which pretty much guarantees that the meeting will be a waste of time). Knowing how to run a meeting efficiently is a core skill for anyone who wants to get things done. Knowing how to take part in someone else's meeting is just as important, but gets far less attention: everyone offers leadership training, nobody offers followership training.

The rules for running meetings quickly and smoothly are well known but rarely followed:

**Decide if there actually needs to be a meeting.** If the only purpose is to share information, have everyone send a brief email instead. Remember, people can read faster than they can speak: if someone has facts for the rest of the team to absorb, the most polite way to communicate them is to type them in.

**Write an agenda.** If nobody cares enough about the meeting to prepare a point-form list of what's to be discussed, the meeting itself probably doesn't need to happen.

**Include timings in the agenda.** Timings help you prevent early items stealing time from later ones. Your first estimates with any new group will be wildly optimistic, so revise them upward for subsequent meetings. However, you shouldn't plan a second or third meeting just because the first one ran over-time: instead, try to figure out why you're running over and fix the underlying problem.

**Prioritize.** Tackle issues that will have high impact but take little time first, and things that will take more time but have less impact later. That way, if the first things run over time, the meeting will still have accomplished something.

**Make one person responsible for keeping things moving.** One person should be made chair, and be responsible for keeping items to time, chiding people who are having side conversations or checking email, and asking people who are talking too much to get to the point. The chair should *not* do all the talking; in fact, whoever is in charge will talk less in a well-run meeting than most other participants.

**Require politeness.** No one gets to be rude, no one gets to ramble, and if someone goes off topic, it's the chair's job to say, "Let's discuss that elsewhere."

**No interruptions.** Participants should raise a finger, put up a sticky note, or make one of the other gestures people use at high-priced auctions when they want to speak. The chair should keep track of who wants to speak and give them time in turn.

**No side conversations.** This makes meetings more efficient because Nobody can actually pay attention to two things at once. If distractions are tolerated, people will miss things or they'll have to be repeated. More importantly, not paying attention is insulting—chatting with a friend while someone explains what they did last week is a really effective way to say, "I don't think you or your work is important."

**No technology.** Insist that everyone put their phones, tablets, and laptops into politeness mode (i.e., close them). If this is too stressful, let participants hang on to their electronic pacifiers but turn off the network so that they really *are* using them just to take notes or check the agenda. The one exception is accessibility needs: if someone needs their phone, their laptop, or some other aid in order to take part in the meeting, nobody has



a right to tell them “no”.

**Take minutes.** Someone other than the chair should take point-form notes about the most important information that was shared, and about every decision that was made or every task that was assigned to someone.

**Take notes.** While other people are talking, participants should take notes of questions they want to ask or points they want to make. (You’ll be surprised how smart it makes you look when it’s your turn to speak.)

**End early.** If your meeting is scheduled for 10:00-11:00, aim to end at 10:55 to give people time to get where they need to go next.

## 20.5 What should we do when the meeting is over?

As soon as the meeting is over, circulate the minutes (i.e., emailed them to everyone or post them to a wiki):

**People who weren’t at the meeting can keep track of what’s going on.**

You and your teammates all have to juggle tasks from other projects or courses, which means that sometimes you won’t be able to make it to meetings. A wiki page, email message, or blog entry is a much more efficient way to catch up than asking a team mate, “Hey, what did I miss?”

**Everyone can check what was actually said or promised.** More than once, one of us has looked over the minutes of a meeting and thought, “Did I say that?” or, “Wait a minute, I didn’t promise to have it ready then!” Accidentally or not, people will often remember things differently; writing it down gives the team a chance to correct mis-recollection, mis-interpretation, or mis-representation, which can save a lot of anguish later on.

**People can be held accountable at subsequent meetings.** There’s no point making lists of questions and action items if you don’t follow up on them later. If you’re using a ticketing system (Chapter 19), create a ticket for each new question or task right after the meeting and update those that are being carried forward. That way, your agenda for the next meeting can start by rattling through the list of open tickets.

Run all your meetings like this for a month, with the goal of making each one a minute shorter than the one before, and we promise that you’ll build better software.

## 20.6 How can we keep people from talking too much or too little?

Some people are so used to the sound of their own voice that they will talk half the time no matter how many other people are in the room. One way to combat this is to give everyone three sticky notes at the start of the meeting. Every time they speak, they have to give up one sticky note. When they're out of stickies, they aren't allowed to speak until everyone has used at least one, at which point everyone gets all of their sticky notes back. This ensures that nobody talks more than three times as often as the quietest person in the meeting, which completely changes group dynamics: people who have given up trying to be heard because they always get trampled suddenly have space to contribute, and the overly-frequent speakers quickly realize just how unfair they have been.

Another useful technique is called interruption bingo. Draw a grid and label the rows and columns with the participants' names. Each time one person interrupts another, add a tally mark to the appropriate cell; halfway through the meeting, take a moment to look at the results. In most cases, you will see that one or two people are doing all of the interrupting, often without being aware of it. After that, saying, "All right, I'm adding another tally to the bingo card," is often enough to get them to throttle back. (Note that this technique is for managing interruptions, not speaking time. It may be completely appropriate for people with more knowledge of a subject to speak about it more often in a meeting, but it is never appropriate to repeatedly cut people off.)

## 20.7 How should we run online meetings?

Chelsea Troy's discussion of why online meetings are often frustrating and unproductive makes an important point: in most online meetings, the first person to speak during a pause gets the floor. As a result, "If you have something you want to say, you have to stop listening to the person currently speaking and instead focus on when they're gonna pause or finish so you can leap into that nanosecond of silence and be the first to utter something. The format...encourages participants who want to contribute to say more and listen less."

The solution is to run a text chat beside the video conference where people can signal that they want to speak. The chair can then select people from the waiting list. This practice can be reinforced by having everyone mute themselves, and only allowing the moderator to unmute people.

## 20.8 What sort of people make teamwork hard?

FIXME: introduction (without Tolstoy)

**Anna** knows more about every subject than everyone else on the team put together—at least, she thinks she does. No matter what you say, she'll correct you; no matter what you know, she knows better. Anna is pretty easy to spot: if you keep track of how often people interrupt one another in team meetings, her score is usually higher than everyone else's put together.

**Bao** is a contrarian: no matter what anyone says, he'll take the opposite side. This can be healthy in small doses, but if someone plays devil's advocate at every turn, they're just impeding progress.

**Caitlin** has so little confidence in her own ability (despite her good grades) that she won't make any decision, no matter how small, until she has checked with someone else. Everything has to be spelled out in detail for her so that there's no possibility of her getting anything wrong.

**Frank** believes that knowledge is power. He enjoys knowing things that other people don't—or to be more accurate, he enjoys it when people know he knows things they don't. Frank is good at making things work, but when asked how he did it, he'll grin and say, "Oh, I'm sure you can figure it out."

**Gary** is a hitchhiker. He has discovered that most people would rather shoulder some extra work than snitch, and he takes advantage of it at every turn. The frustrating thing is that he's so damn *plausible* when someone finally does confront him. "There have been mistakes on all sides," he says, or, "Well, I think you're nit-picking." The only way to deal with Kenny is to stand up to him: remember, if he's not doing his share, *he's the bad guy*, not you.

**Hediyeh** is quiet—very quiet. She never speaks up in meetings, even when she knows that what other people are saying is wrong. She might contribute to the mailing list, but she's very sensitive to criticism, and will always back down rather than defending her point of view. Hediyeh isn't a troublemaker, but rather a lost opportunity.

**Melissa** would easily have made the varsity procrastination team if she'd bothered to show up to tryouts. She means well, and she really does feel bad about letting people down, but somehow something always comes up, and her tasks are never finished until the last possible moment. Of course, that means that everyone who is depending on her can't do their work until *after* the last possible moment...

**Petra** has a favorite phrase: "why don't we". Why don't we write a GUI to help people edit the program's configuration files? Hey, why don't we invent our own little language for designing GUIs? Her energy and enthusiasm are hard to argue with, but argue you must. Otherwise, for every step you move forward, the project's goalposts will recede by two. This is called feature creep, and has ruined many projects that might otherwise have

delivered something small but useful.

**Raj** is rude. “That’s stupid,” and a mixed bag of obscenities are his favorite phrases. “It’s just the way I talk,” he says, not knowing or not caring that for a lot of people, that kind of language has often been a prelude to bullying or worse. His only redeeming grace is that he can’t dissemble as well as Gary, so he’s easier to get rid of.

**Sergei** is simply incompetent. He doesn’t understand the problem, he hasn’t bothered to master the tools and libraries he’s supposed to be using, the code he checks in doesn’t run, and his thirty-second bug fixes introduce more problems than they solve. If he means well, try to re-partition the work so that he’ll do less damage. If he doesn’t, he should be treated like any other hitchhiker.

## 20.9 How should we handle conflict within the team?

FIXME: this is not about harassment or abuse—see Chapter 18 for that.

You just missed an important deadline, and people are unhappy. The sick feeling in the pit of your stomach has turned to anger: you did *your* part, but Marta didn’t finish her stuff until the very last minute, which meant that no one else had time to spot the two big mistakes she’d made. As for Chul, well, he didn’t deliver at all—again. If something doesn’t change, this project is going to pull down your performance review so far that you might have to start looking for a new job.

Situations like this come up all the time. Broadly speaking, there are four ways you can deal with them:

1. Cross your fingers and hope that things will get better on their own, even though the last three times you hoped they would, they didn’t.
2. Do extra to make up for others’ shortcomings. This sometimes works in the short term, and saves you the mental anguish of confronting others, but the time for that “extra” has to come from somewhere; what usually ends up happening is that other parts of the project, or your personal life, suffer.
3. Lose your temper and start shouting. Unfortunately, people often wind up displacing their anger into other parts of their life: I’ve seen developers yell at waitresses for bringing incorrect change when what they really needed to do was tell their boss, “No, I *won’t* work through another holiday weekend to make up for your decision to short-staff the project.”
4. Take constructive steps to fix the underlying problem.

## 20.9. HOW SHOULD WE HANDLE CONFLICT WITHIN THE TEAM?317

Most of us find number four hard because we don't like confrontation. If you manage it properly, though, it is a lot less bruising, which means that you don't have to be as afraid of initiating it. Also, if people believe that you will actually take steps when they bully, lie, procrastinate, or do a half-assed job, they will usually avoid making it necessary. (A colleague once said that she had never met anyone who wished they had waited longer before putting their foot down.)

Here are the steps you should take when you feel that a teammate isn't pulling their weight:

**Make sure you're not guilty of the same sin.** You won't get very far complaining about someone else interrupting in meetings if you do it just as frequently.

**Check expectations.** Are you sure the offender knows what standards they are supposed to be meeting? This is where things like job descriptions or up-front discussion of who's responsible for what come in handy.

**Check the situation.** Is someone dealing with an ailing parent or immigration woes? Have they been put to work on three other projects that you don't know about?

**Document the offense.** Write down what the offender has actually done and why it's not good enough. Doing this will help you clarify matters in your own mind. It's also absolutely necessary if you have to escalate.

**Check with other team members.** Are you alone in feeling that the offender is letting the team down? If so, you aren't necessarily wrong, but it'll be a lot easier to fix things if you have the support of the rest of the team. Finding out who else on the team is unhappy can be the hardest part of the whole process, since you can't even ask the question without letting on that you're upset, and word will almost certainly get back to whoever you're asking about, who might then turn around and accuse you of stirring up trouble. After a couple of unhappy experiences of this kind, I've learned that it's best to raise the issue for the first time in a group of three: you, the person you want to talk to, and a third person who can act as unofficial moderator.

**Talk with the offender.** This should be a team effort: put it on the agenda at a team meeting, present your complaint, and make sure that the offender understands it. In most cases this is enough: human beings are herd animals, and if someone realizes that they're going to be called on their hitchhiking or bad manners, they will usually change their ways.

**Escalate as soon as there's a second offense.** Hitchhikers and others who really don't have good intentions are counting on you giving them one last chance after another until the project is finished and they can go suck the life force out of their next victim. *Don't fall into this trap.* If someone stole your laptop, you'd report it right away. If someone steals your time, you're being pretty generous giving them more than one chance to mend their ways.

In the context of a research project, "escalation" means "taking the issue to your

supervisor”. (If you’re reluctant to do this because you don’t want to be a snitch, go back and re-read the bit about people stealing your laptop.) Of course, your supervisor has probably had dozens of students complain to her over the years about teammates not doing their share—it isn’t uncommon to have both halves of a pair tell the instructor that they’re doing all the work. (This is yet another reason to use version control: it makes it easy to check who’s actually written what.) In order to get her to take you seriously and help you fix your problem, you should send her an email, signed by several people describing the problem and the steps you have already taken to resolve it. Make sure the offender gets a copy as well, and ask your instructor to arrange a meeting to resolve the issue.

This is where documentation is crucial. Hitchhikers are usually very good at appearing reasonable; they’re very likely to nod as you present your case, then say, “Well, yes, but...” and rhyme off a bunch of minor exceptions or cases where others on the team have also fallen short of expectations. If you can’t back up your complaint, your supervisor will likely be left with the impression that the whole team is dysfunctional, and nothing will improve.

One technique your supervisor may ask you to use in a meeting like this is active listening. As soon as one person makes a point, the person on the opposite side of the issue explains it back to them, as in, “So what I think Igor is saying is...” This guarantees that the second person has actually paid attention to what the first person said. It can also defuse a lot of tension, since explaining someone’s position back to them forces you to see the world through their eyes, if only for a few moments.

## 20.10 Summary

- Brown (2007) has lots of good advice on running meetings.
- Brookfield and Preskill (2016) is a catalog of ways to get people talking productively.

## 20.11 Key Points

- FIXME

## Chapter 21

# R Packaging

### 21.1 Questions

- How can I manage the packages my project relies on?
- How can I share R I've written with other people as a package?
- Where and how can I share the code I have created?
- How should I announce that my package exists?
- What extra work do I have to do to meet CRAN's submission criteria?

### 21.2 Objectives

- Create and test a citable, shareable R package.

### 21.3 What's in an R package?

Another response of the wizards, when faced with a new and unique situation, was to look through their libraries to see if it had ever happened before. This was...a good survival trait. It meant that in times of danger you spent the day sitting very quietly in a building with very thick walls.

– Terry Pratchett

The more software you write, the more you realize that a programming language is mostly a way to build and combine software packages. Every widely-used language now has an online repository from which people can download and install packages, and sharing yours is a great way to contribute to the community

that has helped you get where you are. This lesson shows you how to use R's tools to do this.

### 21.3.1 CRAN and Alternatives

CRAN, the Comprehensive R Archive Network, is the best place to find the packages you need. CRAN's famously strict rules ensure that packages run for everyone, but also makes package development a little more onerous than it might be. You can also share packages directly from GitHub, which many people do while packages are still in development. We will explore this in more detail below.

### 21.3.2 Exercise: What packages do you have?

What R packages are currently installed on your computer? How did you figure this out?

### 21.3.3 Exercise: Can you build a package?

1. Clone the GitHub repository for the `here` package at <https://github.com/r-lib/here>.
2. Open `here.Rproj` in RStudio.
3. Build the package.

Does the package build successfully? What messages do you see? Do any of them worry you?

## 21.4 What *is* a package, exactly?

Suppose you have written a useful R script and want to share it with colleagues. You could email it to them (or point them at the GitHub repository it's in), but what should they do next? They could copy the file onto their computer, but then they would have to decide where to put it and remember their decision later on when you updated the script to fix a few bugs. And what about the documentation you so lovingly crafted—where should it go? Oh, and what if different people decide to organize their files differently—how hard will it be to make the things they have built play nicely together?

Packages solve these problems. While the details vary from language to language, packages always require that information about the software be stored in a specific format and in a specific location, so that statements like `library(something)` know where to find what they need. They are rather like



the USB ports of the software world: anything that conforms to a few simple rules can plug in to anything else.

## 21.5 How do I create a package?

An R package must contain the following files:

- The text file **DESCRIPTION** (with no suffix) describes what the package does, who wrote it, and what other packages it requires to run. We will edit its contents as we go along.
- **NAMESPACE**, (whose name also has no suffix) contains the names of everything exported from the package (i.e., everything that is visible to the outside world). As we will see, we should leave its management in the hands of RStudio and the **devtools** package we will meet below.
- Just as **.gitignore** tells Git what files in a project to ignore, **.Rbuildignore** tells R which files to include or not include in the package.
- All of the R source for our package must go in a directory called **R**; sub-directories below this are not allowed.
- As you would expect from its name, the optional **data** directory contains any data we have put in our package. In order for it to be loadable as part of the package, the data must be saved in R's custom **.rda** format. We will see how to do this below.
- Manual pages go in the **man** directory. The bad news is that they have to be in a sort-of-LaTeX format that is only a bit less obscure than the runes inscribed on the ancient dagger your colleague brought back from her latest archeological dig. The good news is, we can embed Markdown comments in our source code and use a tool called **roxygen2** to extract them and translate them into the format that R packages require.
- The **tests** directory holds the package's unit tests. It should contain files with names like `test_some_feature.R`, which should in turn contain functions named `test_something_specific`. We'll have a closer look at these in Chapter ??.

You can type all of this in if you want, but R has a very useful package called **usethis** that will help you create and maintain packages. To show how it works, we will create an R package called **zipffreq** (with no dashes or other special characters in its name) to hold word frequencies in classic English novels. The first step is to load **usethis** in the console with `library(usethis)` and use `usethis::create_package` with the path to the new package directory as an argument:

```
usethis::create_package('~/.zipffreq')
```

```
Setting active project to '/Users/gvwilson/zipffreq'
Creating 'R/'
Creating 'man/'
Writing 'DESCRIPTION'
Writing 'NAMESPACE'
Writing 'zipffreq.Rproj'
Adding '.Rproj.user' to '.gitignore'
Adding '^zipffreq\\.Rproj$', '^\\.Rproj\\.user$' to '.Rbuildignore'
Opening new project 'zipffreq' in RStudio
```

Every well-behaved package should have a README file, a license, and a Code of Conduct, so we will ask `usethis` to add those:

```
usethis::use_readme_md()
usethis::use_mit_license(name="Merely Useful")
usethis::use_code_of_conduct()
```

(Note that `use_mit_license` creates two files: `LICENSE` and `LICENSE.md`. The rules for R packages require the former, but GitHub expects the latter.) We then edit `README.md` to be:

```
# zipffreq
```

An example package for Merely Useful that checks Zipf's Law for classic English novels

```
## Installation
```

```
TBD
```

```
## Example
```

```
FIXME: add an example.
```

and make a similar edit to `DESCRIPTION` so that it contains:

```
Package: zipffreq
Title: Checks Zipf's Law for classic English novels.
Version: 0.0.0.9000
Authors@R:
  person(given = "Merely",
         family = "Useful",
         role = c("aut", "cre"),
         email = "merely.useful@gmail.com")
Description: >
  Contains data and functions for checking Zipf's Law
  for a set of classic English novels.
```

```
License: MIT + file LICENSE
Encoding: UTF-8
LazyData: true
```

We can now go to the **Build** tab in RStudio and run **Check** to see if our empty package makes sense. When we do, the check warns us that there shouldn't be a period at the end of the package title. Once we fix that, we get a clean bill of health.

### 21.5.1 Leftovers

Running **Check** creates another directory called `zipffreq.Rcheck` and a file called `zipffreq_0.0.0.9000.tar.gz`. These are created beside our project directory rather than in it so as not to confuse version control—we don't want the files we are building committed to our repository.

We can now create a placeholder for one of the functions we want to write in a file called `R/frequency.R` either by using **File...New** in RStudio or by running `usethis::use_r('frequency.R')` (which always creates the file in the `R` directory):

```
word_count <- function(word, text) {
  0
}
```

**Build...Check** runs a lot more checks now because we have some actual code for it to look at. It warns us that our function needs documentation, so we will look at that next.

### 21.5.2 Exercise: Packaging the CO2 functions

1. Create a brand-new package to hold the CO2 data and functions for analyzing it.
2. Put the license file where R's packaging rules wants it.
3. Read the documentation for `usethis::use_citation` and create a citation file.
4. Build the package: what warnings do you get and what do they mean?

### 21.5.3 Exercise: Exploring a package

1. Clone a GitHub repository that contains an R package, such as <https://github.com/r-lib/usethis>.
2. Which files in this repository do you recognize, and what are they for?

3. Which files are there to satisfy something other than R's packaging system?

### 21.5.4 Exercise: Ignoring Files

What does `usethis::use_build_ignore` do? When would you use it?

## 21.6 How do I use a package while I'm creating it?

The `devtools` package includes functions to help you create and test packages. After you install it, `devtools::load_all("path/to/package")` will do the same thing that `library(package)` does, but use the files you are developing. (If you run `load_all()` without a path, it will re-load the package in the current working directory, which is what you usually want to do during development.)

`devtools::load_all` is different from `devtools::install`, which not only installs your package, but also tries to install its dependencies from CRAN. After you use `install`, the package will be available to other projects on your computer. This is therefore usually one of the last things you do as you come to the end of development.

## 21.7 How do I document a R package?

Our next task is to document our function. To do this, we turn to Hadley Wickham's *R Packages* and Karl Broman's "R package primer" for advice on writing roxygen2. We then return to our source file and put a specially-formatted comment in front of our code:

```
#' Count how often a word appears in a piece of text
#' 
#' @param word the word to search for
#' @param text the text to search in
#' 
#' @return number of times the word appears in the text
#' 
#' @export

word_count <- function(word, text) {
  0
}
```

This comment text, is an example of embedded documentation, can be easily inserted by placing the cursor somewhere in the function and clicking the button “Code -> Insert Roxygen Skeleton” in RStudio. Over the years, programmers have found that if they put code in one file and documentation in another, the documentation quickly falls out of date with the code: people will change a function’s name or add a new parameter and forget to update the docs. If the documentation lives right beside the code, on the other hand, the next person to modify the code is far more likely to remember to update it. Tools like roxygen2 can read the code, extract the documentation, and format it as HTML or PDF. They can also do things like create an index, which would be even more painful to do by hand than writing the documentation itself.

roxygen2 processes comment lines that start with `#'` (hash followed by single quote). Putting a comment block right before a function associates that documentation with that function, and `@something` indicates a roxygen2 command, so what this file is saying is:

- the function has two parameters called `word` and `text`
- it returns the number of times the word is found in the text; and
- we want to export it (i.e., we want it to be visible outside the package).

These roxygen2 text can also be written with Markdown formatting, so that for instance using `**word**` will bold the word. Other common formatting includes:

- `[function_name()]` or `[object_name]` to link to other function or object, respectively, documentation in the same package.
- `[pkgname::function_name()]` or `[pkgname::object_name]` to link to function or object, respectively, documentation in another package.
- `[link name](website-link)` to link to a website.

To enable Markdown with roxygen2 we’ll need to install the package roxygen2md package and use:

```
roxygen2md::roxygen2md("full")
```

which will output something like:

```
Setting active project to '/Users/gvwilson/zipffreq'
No files changed.
0 source files changed
Running `devtools::document()`
Updating zipffreq documentation
Writing NAMESPACE
Loading zipffreq
Writing NAMESPACE
Review the changes carefully
Commit the changes to version control
Setting active project to '<no active project>'
```

This converts all existing documentation to use Markdown version and will also add this line to the DESCRIPTION file:

```
Roxygen: list(markdown = TRUE)
```

Ok, our function is now documented, but when we run `Check`, we still get a warning. After a bit more searching and experimentation, we discover that we need to run `devtools::document()` to regenerate documentation because it isn't done automatically. When we do this, we get:

```
Updating zipffreq documentation
Updating roxygen version in /Users/gvwilson/zipffreq/DESCRIPTION
Writing NAMESPACE
Loading zipffreq
Writing NAMESPACE
Writing word_count.Rd
```

We now have two files to look at: `NAMESPACE` and `word_count.Rd`. The first looks like this:

```
# Generated by roxygen2: do not edit by hand

export(word_count)
```

The comment at the start tells roxygen2 it can overwrite the file, and reminds us that we shouldn't edit it by hand. The `export(word_count)` directive is what we really want: as you might guess from the name `export`, it tells the package builder to make this function visible outside the package.

What about `word_count.Rd`? It lives in the `man` directory and now contains:

```
% Generated by roxygen2: do not edit by hand
% Please edit documentation in R/frequency.R
\name{word_count}
\alias{word_count}
\title{Count how often a word appears in a piece of text}
\usage{
word_count(word, text)
}
\arguments{
\item{word}{the word to search for}

\item{text}{the text to search in}
}
\value{
number of times the word appears in the text
}
\description{
Count how often a word appears in a piece of text
```

```
}
```

Again, there's a comment at the top to remind us that we shouldn't edit this by hand. After doing this, we go into "Build... More... Configure build tools" and check "Generate documentation with Roxygen". Running **Check** again now gives us a clean bill of health. If we use the **Install and Restart** button in RStudio's **Build** tab, we can now use `?word_count` in the console to view our help.

### 21.7.1 Exercise: Document a function

FIXME: write this exercise once the CO2 manipulation functions exist.

### 21.7.2 Exercise: Use Markdown documentation

Read the documentation for the `roxygen2md` package, which allows you to write R documentation in Markdown, and then convert `zipffreq` to use Markdown documentation.

## 21.8 What should I document?

The answer to the question in this section's title depends on what stage of development you are in. If you are doing exploratory programming, a one-line comment to remind yourself of each function's purpose is good enough. (In fact, it's probably better than what most people do.) That comment should begin with an active verb and describe how inputs are turned into outputs. If the function has any side effects, you should eliminate them. If you can't, you should describe them too.

An active verb is something like "extract", "normalize", or "find". For example, these are all good one-line comments:

- "Create a list of current ages from a list of birth dates."
- "Ensure training parameters lie in [0..1]."
- "Reduce the red component of each pixel."

You can tell your one-liners are useful if you can read them aloud in the order the functions are called in place of the function's name and parameters.

Once you start writing code for other people—including yourself three months from now—your documentation should describe:

1. The name and purpose of every function and constant in your code.
2. The name, purpose, and default value (if any) of every parameter to every function.

3. Any side effects the function has.
4. The type of value returned by every function.
5. Why and how the function will deliberately fail. If a function uses something like `stopifnot` or `assert` to check that a condition holds, then that halt-and-catch-fire behavior is effectively part of its interface.

### 21.8.1 Exercise: Fixing documentation

FIXME: provide poorly-documented function in CO2 package and ask learner to find and fix gaps.

### 21.8.2 Exercise: Cross-references

1. What should you add to the roxygen2 comments for one function to link to the documentation for another function?
2. Add a cross-reference from the documentation for the function FIXME in the CO2 package to the documentation for the function FIXME.

### 21.8.3 Exercise: Documenting error conditions

The guidelines above said that authors should document why and how their functions will deliberately fail. Where and how should you do this using roxygen2 for R?

## 21.9 How do I manage package dependencies?

In order to understand the rest of what follows, it's important to understand that R packages are distributed as compiled byte code, *not* as source code (which is how Python does it) or as binary code (which is how app stores distribute things). When a package is built, R loads and checks the code, then saves the corresponding low-level instructions. Our R files should therefore define functions, not run commands immediately: if they do the latter, those commands will be executed every time the script loads, which is probably not what users will want.

As a side effect, this means that if a package uses `load(something)`, then that `load` command is executed *while the package is being compiled*, and *not* while the compiled package is being loaded by a user after distribution (Figure 21.1). If we have loaded the library by hand in our R session during development, though, we might not notice the problem.

How then can our packages use other packages? The safest way is to use fully-qualified names such as `stringr::str_replace` every time we call a function





Figure 21.1: Package Distribution

from another package or that is defined somewhere outside our package. Let's modify our word counter to use `stringr::str_count`:

```
word_count <- function(word, text) {
  stringr::str_count(text, word)
}
```

Since our compiled-and-distributable package will only contain the compiled code for its own functions, direct calls to functions from other packages won't work after the package is installed. To fix this, we ask `usethis` to add a note about `stringr` to `DESCRIPTION`:

```
usethis::use_package('stringr')
```

The bottom of `DESCRIPTION` now has these two lines:

```
Imports:
  stringr
```

The `Imports` field in `DESCRIPTION` tells R what else it has to install when it installs our package. To be explicit about the version of the package, we can run:

```
usethis::use_tidy_versions()
```

Which will add the version of the package to the `DESCRIPTION` file, which will look like:

```
Imports:
  stringr (>= 1.4.0)
```

(See Section 12.9 for a discussion of version numbering.)

All right: are we done now? No, we are not:

### 21.9.1 Exercise: Document dependencies

1. Modify the functions in the CO2 package to use `package::function name` for everything.
2. Modify the DESCRIPTION file to document the package's dependencies.

### 21.9.2 Importing

1. What does `@import` do in roxygen2 documentation?
2. When should or shouldn't you use it?

## 21.10 How can I share my package via GitHub?

We said in Section 21.3 that R packages could be shared through GitHub as well as through CRAN. If someone has done this, installing the package on your computer is as simple as:

```
devtools::install_github("username/reponame")
```

where `username` and `reponame` are the names of the user and the project respectively. If you want to share your work, all you have to do is create a repository whose contents are laid out as described earlier in this chapter: if it looks like an R project, `install_github` will treat it as one.

## 21.11 How can I add data to a package?

The last steps are to add some cleaned-up data to our package and document the package as a whole. If we want to create fake data or clean up raw data, we should first run the command:

```
# Let's name the data "small_data", but it can be anything.
usethis::use_data_raw("small_data")
```

which outputs:

```
Creating 'data-raw/'
Adding '^data-raw$' to '.Rbuildignore'
Writing 'data-raw/small_data.R'
Modify 'data-raw/small_data.R'
```

Finish the data preparation script in 'data-raw/small\_data.R'

Use ``usethis::use_data()`` to add prepared data to package

that creates the `data-raw/` folder and the `small_data.R` file in that folder.  
Inside that new file is:

```
## code to prepare `small_data` dataset goes here
```

```
usethis::use_data("small_data")
```

Let's replace the first comment with:

```
small_data <- tribble(
  ~word, ~ count,
  "some", 2,
  "words", 1,
  "appear", 1,
  "times", 1)
```

which looks like:

```
# A tibble: 4 x 2
  word    count
  <chr>   <dbl>
1 some         2
2 words        1
3 appear       1
4 times        1
```

Remove the quotes within the `use_data()`. The new file should contain the code:

```
small_data <- tribble(
  ~word, ~ count,
  "some", 2,
  "words", 1,
  "appear", 1,
  "times", 1
)

usethis::use_data(small_data)
```

Run the script by hitting "Source" (or `source("data-raw/small_data.R")`).  
This will output:

```
Creating 'data/'
Saving 'small_data' to 'data/small_data.rda'
```

Now the `small_data` data.frame object is saved to `data/` and is now accessible throughout the package!

When we run `Check`, we get a complaint about an undocumented data set, so we create a file called `R/small_data.R` to hold documentation about the dataset and put this in it:

```
#' Sample word frequency data.
#'  
#'This small dataset contains word frequencies for tutorial purposes.  
#'  
#'@docType data  
#'  
#'@format A data frame  
#'\describe{  
#'  \item{word}{The word being counted (chr)}  
#'  \item{n}{The number of occurrences (int)}  
#'}  
"small_data"
```

Everything except the last line is a roxygen2 comment block that describes the data in plain language, then uses some tags and directives to document its format and fields. The line `@docType` tells roxygen2 that this comment describes data rather than a function, and the last line is the string `"small_data"`, i.e., the name of the dataset. We will create one placeholder R file like this for each of our datasets, and each will have that dataset's name as the thing being documented.

Running `Check` now gives us a different warning:

```
Warning: package needs dependence on R (>= 2.10)
```

We do *not* fix this by adding another line under `Imports` in `DESCRIPTION`, since R itself isn't a package. Instead, we add this line:

```
Depends: R (>= 2.10)
```

and get a clean bill of health.

We use a similar trick to document the package as a whole: we create a file `R/zipffreq.R` (i.e., a file with exactly the same name as the package) and put this in it:

```
#' Example of an R package.
#'  
#'@author Merely Useful, \email{merely.useful@gmail.com}  
#'@docType package  
#'@name zipffreq  
NULL
```

That's right: to document the entire package, we document `NULL`. One last build, and our package is ready to deliver.

### 21.11.1 The Virtues of Laziness

We should always put `LazyData: TRUE` in `DESCRIPTION` so that datasets are only loaded on demand.

### 21.11.2 Exercise: Letting `usethis` do even more work

What does `usethis::use_package_doc` do?

### 21.11.3 Exercise: Scripting data creation

1. Put the code to create `small_data` in an R script in `data-raw` via `usethis::use_data_raw()`.
2. Include `usethis::use_data(small_data)` at the end of the script.
3. Describe what this script does step by step.

### 21.11.4 Exercise: Add sample data to the `CO2` package

1. Create a small sample of the `CO2` data in a tibble called `sample_CO2_data`.
2. Save it and document it.

### 21.11.5 Exercise: Reproducible data set creation

1. Write a small script to create and save the sample dataset from the previous exercise.
2. Where should this script go in your package? How should you document its existence and usage?
3. Swap packages with a colleague. Can you regenerate their sample data using only what is (documented) in their package?

## 21.12 Summary

## 21.13 Key Points

- Packages allow software to be shared in manageable ways.
- R packages can be shared through CRAN or GitHub, or managed locally during development.
- Packages can contain code and data.
- A package must contain `DESCRIPTION` and `NAMESPACE` files.
- Use `.Rbuildignore` to control what is and isn't included in a package.

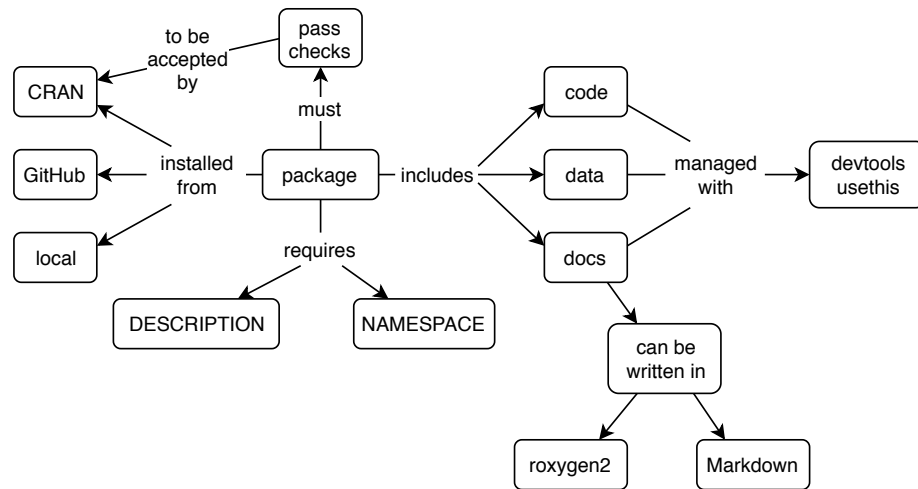


Figure 21.2: R Packaging Concept Map

- Include a README, a license, and a citation file in every package.
- Use `usethis` and `devtools` to manage package development.
- Put documentation in the `man` directory and tests in the `tests` directory.
- Use `roxygen2` or `Markdown` to document the contents of a package.

## Chapter 22

# Python Packaging

### 22.1 Questions

- How can I manage the libraries my project relies on?
- How can I package up my work for others to use?
- How should I announce my work?

### 22.2 Objectives

- Create and use virtual environments to manage library versions without conflict.
- Create and test a citable, shareable Pip package.

### 22.3 Introduction

Another response of the wizards, when faced with a new and unique situation, was to look through their libraries to see if it had ever happened before. This was...a good survival trait. It meant that in times of danger you spent the day sitting very quietly in a building with very thick walls.

– Terry Pratchett

The more software you write, the more you realize that a programming language is a way to build and combine software libraries. Every widely-used language now has an online repository from which people can download and install those

libraries. This lesson shows you how to use Python's tools to create and share libraries of your own.

This material is based in part on Python 102 by Ashwin Srinath.

## 22.4 How can I turn a set of Python source files into a module?

Any Python source file can be imported by any other. When a file is imported, the statements in it are executed as it loads. Variables defined in the file are then available as `module.name`. (This is why Python files should be named using `pothole_case` instead of `kebab-case`: an expression like `import some-thing` isn't allowed because `some-thing` isn't a legal variable name.)

As an example, we can put a constant and two functions used in our Zipf's Law study in a file called `zipf.py`:

```
from pytest import approx

RELATIVE_ERROR = 0.05

def make_zipf(length):
    assert length > 0, 'Zipf distribution must have at least one element'
    result = [1/(1 + i) for i in range(length)]
    return result

def is_zipf(hist, rel=RELATIVE_ERROR):
    assert len(hist) > 0, 'Cannot test Zipfiness without data'
    scaled = [h/hist[0] for h in hist]
    perfect = make_zipf(len(hist))
    return scaled == approx(perfect, rel=rel)
```

and then use `import zipf`, `from zipf import is_zipf`, and so on:

```
from zipf import make_zipf, is_zipf

generated = make_zipf(5)
print('generated distribution: {}'.format(generated))
generated[-1] *= 2
print('passes test with default tolerance: {}'.format(is_zipf(generated)))
print('passes test with tolerance of 1.0: {}'.format(is_zipf(generated, rel=1.0)))
```

Running this program produces the following output:



```
generated distribution: [1.0, 0.5, 0.3333333333333333, 0.25, 0.2]
passes test with default tolerance: False
passes test with tolerance of 1.0: True
```

It also creates a sub-directory called `__pycache__` that holds the compiled versions of the imported files. The next time Python imports `zipf`, it checks the timestamp on `zipf.py` and the timestamp on the corresponding file in `__pycache__`. If the latter is more recent, Python doesn't bother to recompile the file: it just loads the bytes in the cached version and uses those. To avoid confusing it, we (almost) always put `__pycache__` in `.gitignore`.

## 22.5 How can I control what is executed during import and what isn't?

Sometimes it's handy to be able to import code and also run it as a program. For example, we may have a file full of useful functions for extracting keywords from text that we want to be able to use in other programs, but also want to be able to run `keywords somefile.txt` to get a listing.

To help us do this (and other things we'll see later), Python automatically creates a variable called `__name__` in each module. If the module is the main program, that variable is assigned the string `'__main__'`. Otherwise, it is assigned the module's name. Using this leads to modules like this:

```
import sys
from pytest import approx

USAGE = '''zipf num [num...]: are the given values Zipfy?'''
RELATIVE_ERROR = 0.05

def make_zipf(length):
    ...as before...

def is_zipf(hist, rel=RELATIVE_ERROR):
    ...as before...

if __name__ == '__main__':
    if len(sys.argv) == 1:
        print(USAGE)
    else:
        values = [int(a) for a in sys.argv[1:]]
```

```
result = is_zipf(values)
print('{}: {}'.format(result, values))
sys.exit(0)
```

Here, the code guarded by `if __name__ == '__main__':` isn't executed when the file is loaded by something else. We can test this by re-running `use.py` as before: the usage message doesn't appear, which means the main block wasn't executed, which is what we want.

## 22.6 How can I install a Python package?

The most common way to install Python packages is to use a tool called `pip`. The command `pip install package` checks to see if the package is already installed (or needs to be upgraded); if so, it downloads the package from PyPI (the Python Package Index), unpacks it, and installs it. Depending on where Python is installed, completing that installation may require administrative privileges; for example, if Python is installed in `/usr/bin/python`, you may need to run `sudo` to overwrite previously-installed libraries. *This is a bad idea*, since system tools may depend on particular versions of those packages, and may break if you overwrite them. Section 22.10 shows how to avoid these problems.

Since a project may depend on many packages, developers frequently put a list of those dependencies in a file called `requirements.txt`. `pip install -r requirements.txt` will then install the dependencies listed in that file. (The file can be called anything, but everyone uses `requirements.txt`, so you should too.) This file can just list package names, or it can specify exact versions, minimum versions, etc.:

```
request
scipy==1.1.0
tdda>=1.0
```

If you want to create a file like this, `pip freeze` will print the exact versions of all installed packages. You usually won't use this directly in `requirements.txt`, since your project probably doesn't depend on all of the listed files, but it's a good practice to save this in version control when producing reports so that you can reproduce your results later (Chapter 24).

## 22.7 How can I create an installable Python package?

Packages have to come from somewhere, and that “somewhere” is mostly developers like you. Creating a Python package is fairly straightforward, and mostly

comes down to having the right directory structure.

A package is a directory that contains a file called `__init__.py`, and may contain other files or sub-directories containing files. `__init__.py` can contain useful code or be empty, but either way, it has to be there to tell Python that the directory is a package. For our Zipf example, we can reorganize our files as follows:

```
+-- use.py
+-- zipf
    +- __init__.py
```

`zipf/__init__.py` contains `RELATIVE_ERROR` and the functions we've seen before. `use.py` doesn't change—in particular, it can still use `import zipf` or `from zipf import is_zipf` even though there isn't a file called `zipf.py` any longer. When we run `use.py`, it loads `zipf/__init__.py` when `import zipf` executes and creates a `__pycache__` directory inside `zipf`.

## 22.8 How can I manage the source code of large packages?

As our package grows, we should split its source code into multiple files. To show how this works, we can put the Zipf generator in `zipf/generate.py`:

```
def make_zipf(length):
    assert length > 0, 'Zipf distribution must have at least one element'
    result = [1/(1 + i) for i in range(length)]
    return result
```

and then import that file in `__init__.py`.

```
import sys
from pytest import approx
from . import generate

RELATIVE_ERROR = 0.05

def is_zipf(hist, rel=RELATIVE_ERROR):
    assert len(hist) > 0, 'Cannot test Zipfiness without data'
    scaled = [h/hist[0] for h in hist]
    perfect = generate.make_zipf(len(hist))
    return scaled == approx(perfect, rel=rel)
```

We write that import as `from . import generate` to make sure that we will get the `generate.py` file in the same directory as `__init__.py` (the `.` means

“current directory”, just as it does in the Unix shell).

When we arrange our code like this, the code that *uses* the library must refer to `zipf.generate.make_zipf`:

```
import zipf

generated = zipf.generate.make_zipf(5)
print('generated distribution: {}'.format(generated))
generated[-1] *= 2
print('passes test with default tolerance: {}'.format(zipf.is_zipf(generated)))
print('passes test with tolerance of 1.0: {}'.format(zipf.is_zipf(generated, rel=1.0)))
```

## 22.9 How can I distribute software packages that I have created?

People can always get your package by cloning your repository and copying files from that (assuming your repository is accessible, which is should be for published research), but it’s much friendlier to create something they can install. For historical reasons, Python has several ways to do this. We will show how to use `setuptools`, which is the lowest common denominator; `conda` is a modern does-everything solution, but has larger startup overhead.

To use `setuptools`, we must create a file called `setup.py` in the directory *above* the root directory of the package:

```
+-- setup.py
+- use.py
+- zipf
  +- __init__.py
  +- generate.py
```

The file `setup.py` must have exactly that name, and must contain these lines:

```
from setuptools import setup, find_packages

setup(
    name='zipf',
    version='0.1',
    author='Greg Wilson',
    packages=find_packages()
)
```

The `name`, `version`, and `author` parameters to `setup` are self-explanatory; you should modify the values assigned to them for your package; the function `find_packages` returns a list of things worth packaging.

## 22.9. HOW CAN I DISTRIBUTE SOFTWARE PACKAGES THAT I HAVE CREATED?341

Once you have created this file, you can run `python setup.py sdist` to create your package. The verb `sdist` stands for “source distribution”, meaning that the source of the Python files is included in the package:

```
$ python setup.py sdist

running sdist
running egg_info
creating zipf.egg-info
writing zipf.egg-info/PKG-INFO
writing dependency_links to zipf.egg-info/dependency_links.txt
writing top-level names to zipf.egg-info/top_level.txt
writing manifest file 'zipf.egg-info/SOURCES.txt'
reading manifest file 'zipf.egg-info/SOURCES.txt'
writing manifest file 'zipf.egg-info/SOURCES.txt'
warning: sdist: standard file not found: should have one of README, README.rst, README.txt, README

running check
warning: check: missing required meta-data: url

warning: check: missing meta-data: if 'author' supplied, 'author_email' must be supplied too

creating zipf-0.1
creating zipf-0.1/zipf
creating zipf-0.1/zipf.egg-info
copying files to zipf-0.1...
copying setup.py -> zipf-0.1
copying zipf/__init__.py -> zipf-0.1/zipf
copying zipf/generate.py -> zipf-0.1/zipf
copying zipf.egg-info/PKG-INFO -> zipf-0.1/zipf.egg-info
copying zipf.egg-info/SOURCES.txt -> zipf-0.1/zipf.egg-info
copying zipf.egg-info/dependency_links.txt -> zipf-0.1/zipf.egg-info
copying zipf.egg-info/top_level.txt -> zipf-0.1/zipf.egg-info
Writing zipf-0.1/setup.cfg
creating dist
Creating tar archive
removing 'zipf-0.1' (and everything under it)
```

We will look at how to clean up the warnings about `README.md`, `url`, and `author_email` in the exercises.

`python setup.py sdist` creates a compressed file `dist/zipf-0.1.tar.gz` that contains the following:

```
$ tar ztvf dist/zipf-0.1.tar.gz

drwxr-xr-x  0 ptery staff      0 20 Aug 15:36 zipf-0.1/
-rw-r--r--  0 ptery staff    180 20 Aug 15:36 zipf-0.1/PKG-INFO
-rw-r--r--  0 ptery staff     38 20 Aug 15:36 zipf-0.1/setup.cfg
```

```

-rw-r--r-- 0 ptery staff 145 20 Aug 13:40 zipf-0.1/setup.py
drwxr-xr-x 0 ptery staff 0 20 Aug 15:36 zipf-0.1/zipf/
-rw-r--r-- 0 ptery staff 317 20 Aug 13:34 zipf-0.1/zipf/__init__.py
-rw-r--r-- 0 ptery staff 163 20 Aug 13:34 zipf-0.1/zipf/generate.py
drwxr-xr-x 0 ptery staff 0 20 Aug 15:36 zipf-0.1/zipf.egg-info/
-rw-r--r-- 0 ptery staff 1 20 Aug 15:36 zipf-0.1/zipf.egg-info/dependency_links.txt
-rw-r--r-- 0 ptery staff 180 20 Aug 15:36 zipf-0.1/zipf.egg-info/PKG-INFO
-rw-r--r-- 0 ptery staff 154 20 Aug 15:36 zipf-0.1/zipf.egg-info/SOURCES.txt
-rw-r--r-- 0 ptery staff 5 20 Aug 15:36 zipf-0.1/zipf.egg-info/top_level.txt

```

The source files `__init__.py` and `generate.py` are in there, along with the odds and ends that `pip` will need to install this package properly when the time comes.

## 22.10 How can I manage the packages my projects need?

We want to test the package we just created, but we *don't* want to affect the packages we already have installed, and as noted earlier, we may not have permission to write into the directory that contains system-wide packages. (For example, we may be testing something out on a cluster shared by our whole department.) The solution is to use a virtual environment. These are slowly being superceded by more general solutions like Docker, but they are still the easiest solution for most of us.

A virtual environment is a layer on top of an existing Python installation. Whenever Python needs to find a library, it looks in the virtual environment first. If it can satisfy its needs there, it's done; otherwise, it looks in the underlying environment. This gives us a place to install packages that only some projects need, or that are still under development, without affecting the main installation.

FIXME: figure

We can create and manage virtual environments using a tool called `virtualenv`. To install it, run `pip install virtualenv`. Once we have done that, we can create a new virtual environment called `test` by running:

```
$ virtualenv test
```

```

Using base prefix '/Users/pterry/anaconda3'
New python executable in /Users/pterry/test/bin/python
Installing setuptools, pip, wheel...
done.

```

`virtualenv` creates a new directory called `test`, which contains sub-directories called `bin`, `lib`, and so on—everything needed for a minimal Python installation.

## 22.10. HOW CAN I MANAGE THE PACKAGES MY PROJECTS NEED?343

Crucially, `test/bin/python` checks for packages in `test/lib` *before* checking the system-wide install.

We can switch to the `test` environment by running:

```
$ source test/bin/activate
```

`source` is a Unix shell command meaning “run all the commands from a file in this currently-active shell”. We use it because typing `test/bin/activate` on its own would run those commands in a sub-shell, which would have no effect on the shell we’re in. Once we have done this, we’re running the Python interpreter in `test/bin`:

```
$ which python
```

```
/Users/pterry/test/bin/python
```

We can now install packages to our heart’s delight. Everything we install will go under `test`, and won’t affect the underlying Python installation. When we’re done, we can switch back to the default environment with `deactivate`. (We don’t need to `source` this.)

Most developers create a directory called `~/envs` (i.e., a directory called `envs` directly below their home directory) to store their virtual environments:

```
$ cd ~
```

```
$ mkdir envs
```

```
$ which python
```

```
/Users/pterry/anaconda3/bin/python
```

```
$ virtualenv envs/test
```

```
Using base prefix '/Users/pterry/anaconda3'
```

```
New python executable in /Users/pterry/envs/test/bin/python
```

```
Installing setuptools, pip, wheel...done.
```

```
$ which python
```

```
/Users/pterry/anaconda3/bin/python
```

```
$ source envs/test/bin/activate
```

```
(test)
```

```
$ which python
```

```
/Users/pterry/envs/test/bin/python
```

```
(test)
```

```
$ deactivate
```

```
$ which python
```

```
/Users/pterry/anaconda3/bin/python
```

Notice how every command now displays `(test)` when that virtual environment is active. Between Git branches and virtual environments, it can be very easy to lose track of what exactly you're working on and with. Having prompts like this can make it a little less confusing; using virtual environment names that match the names of your projects (and branches, if you're testing different environments on different branches) quickly becomes essential.

## 22.11 How can I test package installation?

Now that we have a virtual environment set up, we can test the installation of our Zipf package:

```
$ pip install ./src/package/05/dist/zipf-0.1.tar.gz
Processing ./src/package/05/dist/zipf-0.1.tar.gz
Building wheels for collected packages: zipf
  Running setup.py bdist_wheel for zipf ... done
  Stored in directory: /Users/pterry/Library/Caches/pip/wheels/6b/de/80/d72bb0d6e7c65b
Successfully built zipf
Installing collected packages: zipf
Successfully installed zipf-0.1
(test)

$ python
>>> import zipf
>>> zipf.RELATIVE_ERROR
0.05

$ pip uninstall zipf
Uninstalling zipf-0.1:
  Would remove:
    /Users/pterry/envs/test/lib/python3.6/site-packages/zipf-0.1.dist-info/*
    /Users/pterry/envs/test/lib/python3.6/site-packages/zipf/*
Proceed (y/n)? y
  Successfully uninstalled zipf-0.1
(test)
```

Again, one environment per project and one project per environment might use more disk space than is absolutely necessary, but it will still be less than most of your data sets, and will save a lot of debugging.

## 22.12 Announcing Work

FIXME: <https://medium.com/indeed-engineering/marketing-for-data-science-a-7-step-go-to-market-pl>



## 22.13 Summary

FIXME: create concept map for packages

## 22.14 Exercises

### 22.14.1 Clean up warning messages

- FIXME: clean up warning messages from `python setup.py sdist`

## 22.15 Key Points

- Use `virtualenv` to create a separate virtual environment for each project.
- Use `pip` to create a distributable package containing your project's software, documentation, and data.



## Chapter 23

# Continuous Integration

### 23.1 Questions

### 23.2 Objectives

### 23.3 Introduction

Continuous integration (CI) in software development is a simple idea: to frequently merge in (“integrate”) additions of or modifications to code so that a software project gets developed and tested in small, regular increments. The purpose of CI is to catch bugs or problems early before they become bigger problems later on. There are automated tools available that help with building and testing any change to the code in a repository. What the tools do is download the software repository, run and build the software and code on the repository, run any tests, and then print out the results of these builds and tests. That way, you can check if there are any problems with your software and fix any that occur.

The most widely used CI system is Travis CI. It integrates well with Github, and will run tests on multiple platforms and with multiple versions of tools. At a minimum, Travis will build the software. If you want to confirm that your software works as intended, you should create tests as detailed in Chapter ???. If tests exist, Travis will run them and print the results for you to check. Remember, the tests that the CI runs are only as good as you’ve written them.

TODO: Confirm unit test chapter ref.

So why should you use CI? There are several reasons:

1. CI will construct a clean, separate computing environment before building and testing your software. That gives you a better idea on whether potential users of your software can actually use it. Often your software will build and work on *your* computer, but not others' computers because of differences in software versions, operating systems, and other potential dependency conflicts. So part of good practice for creating software is to confirm that it works at least on a clean environment.
2. CI will help you catch problems or bugs sooner, which should reduce your stress and frustration compared to if you have to hunt for a bug later in the software development.
3. CI is particularly useful when working on a project with others, as bugs that arise from someone else's code contribution can be harder for you to resolve. Preventing them from being merged in the first place saves time. By encouraging frequent well-tested merges, CI also reduces the problems that can occur when multiple people wait to integrate their work right before a big milestone and then want to merge in many big changes all at once.
4. CI can also be extended and customized to be used in other situations outside of software projects, like when writing a book or creating a website. Having the knowledge and skill for using CI can allow you to apply these tools to other areas of work and ultimately reduce more manual and tedious work tasks.

In this lesson we will show you how to set up Travis for your project.

## 23.4 The basics of setting up CI on a repository?

The general steps for setting up CI are:

1. Go to the Travis website and get Travis to start watching the repository.
2. Create a `.travis.yml` file in your repository.
3. Add some standard commands to the file, like `language` (e.g. Python or R).
4. Add and commit the `.travis.yml` file into your Git repository.
5. Push to the repository to GitHub.

And that is the basic steps involved! Let's go through it in more detail. The very first thing to do is create a Travis account by linking your GitHub account. Go to Travis and create your profile.

### 23.4.1 Travis is watching (your repo)

Note: This process is for general Travis usage. When working with R projects or packages, this process is simplified by using the

`usethis::use_travis()` command, which we will discuss at the end of this session.

After the “step-zero” of setting up your Travis account, the first step to using Travis is to tell it about your repository. Go to Travis-CI and on the left, besides the “My Repositories”, click the “+” to add a repository that is already on GitHub.

Then, find your repository in the list and flick the switch button on (so it is green). You may have to re-sync the list of repositories with GitHub if you can’t find the repository (green “Sync account” button on the left sidebar).

You’ve now told Travis to watch the repository for any updated commits and run any commands as needed.

### 23.4.2 Configure your project for Travis

Next we need to create a file called `.travis.yml` in the root directory of the repository. The leading `.` in the name is used by Unix-like systems (Mac OS or Linux) to hide the file, but on Windows systems it won’t be hidden. This file contains YAML key-value settings that Travis will use and interpret. A simple template `.travis.yml` configuration file may look like:

```
# Which programming language to use, e.g. Python
language: ...

# Python specific. Sets the Python version
python:
- ...

# Commands to use to install, such as extra software packages
install:
- ...

# Commands to run in order to build or test your software
script:
- ...
```

Let’s break it down what each key does:

- **language** tells Travis which programming language you’re using.
- **python** is specific to when **language** is set to **python** and sets the version you want to use, e.g. “3.6”. You can ask Travis to test our project with several different versions by adding more `-` items.
- **install** tells Travis how to install the software you need for your package. For Python packages, the convention is to put a list of packages in a file called `requirements.txt` for `pip` to use (Chapter 22).

- **script** tells Travis how to actually run the tests. We can put almost anything here, provided the code doesn't need human interaction (no question prompts for answers).

An example file that has the keys filled out might be something like:

```
language: python

python:
- "3.6"

install:
- pip install -r requirements.txt

script:
# if an example build.py file was in the src folder
- python src/build.py
```

After we've added and committed this file to Git, pushed up to GitHub, and (if Travis is watching the GitHub repository) Travis will start. Every time commits are pushed to the GitHub repository, Travis will:

1. Create a new Linux image.
2. Install the desired version of Python.
3. Install the software described in `requirements.txt`.
4. Run the commands in the **script** key, in this case to run `python src/build.py`.
5. Report the results at <https://travis-ci.org/USER/REPO> (USER is your username, and REPO is the name of the repository).

Travis' summary report tells us whether the build passed or failed. It fails when it encounters warnings or errors. In the image below, the main thing to look for is whether the colour is green and says "passed" (as shown), or if the colours are red and the message is "failed" or "error".

The log below this overview contains a ton of information, most of which are only relevant when you start needing to debug why and how your software doesn't build in a clean, non-local environment (aka. not your own computer). We'll get into more detail about interpreting the logs and overviews later in the chapter.

### 23.4.3 Display CI build status in README

Travis's dashboard is very useful, but it would also be nice to have a short, quick display of the Travis build status on our software project's README since that is where people first go to look for information on the project. To add this status "badge", click the "build icon" as shown in the top right corner of Figure

... This will bring up a dialog box where you can select the item “Markdown” in the “Format” menu list. In the “Result” box will be the markdown text you can copy and paste into your repository’s `README.md` file. It’s best to paste the text right below the first `#` header. Commit and push to GitHub and your project will now show the “badge” for the Travis CI build, like shown in the image below.

TODO: Add ref to image?

### 23.4.4 Setting up Travis in R

In R, setting up and using Travis is greatly simplified by using:

```
usethis::use_travis()
```

This command does several things at once. It:

- Creates a typical `.travis.yml` configuration file used for R packages.
- Automatically adds a Travis build badge to your `README.md`.
- Opens up the Travis CI website to the page of your package so you can activate Travis for the repository. Note, this assumes you already have your package on GitHub and that you have a Travis account.

See the `usethis` webpage for more details and help on this function.

You’ve now set up CI for your R or Python software repository! Each new push to your GitHub repository will trigger Travis to run the specified commands. You can now start using it to find errors or problems and to making your software usable to others.

## 23.5 Using Travis to test your software

The primary reason to use CI is to run the package’s tests (as we covered in the [Unit Test] section) and check how those tests perform in a clean environment. So getting Travis to run the unit tests in your package is fairly easy. For R packages, running tests is built into Travis CI so you don’t need to modify anything. For Python packages, the `script` key in the `.travis.yml` file needs to have `pytest` added, like so:

TODO: Reference to unit test chapter.

```
language: python
python:
- "3.6"
install:
- pip install -r requirements.txt
```

```
script:
- pytest
```

After committing and pushing to your GitHub repository, Travis will now always run your tests in its clean environment. If the tests fail, Travis will update the status to “failed” or “error”. Then you can go through the Travis log and start debugging what went wrong, which can sometimes be a very difficult task and is always different depending on the project, tests, and programming language. All details about the build are contained in the Travis “Job log”, right below the build overview. Searching for the problem can sometimes be quite tedious. In the final exercise we’ll get you to intentionally fail the Travis build to see how it looks.

TODO: Do we want to add more explanation about this? I feel running unit tests are basically the biggest reason to use CI.

## 23.6 How can I use CI for non-testing purposes?

This section is meant only as a brief resource and introduction to using CI for more than just testing software. For instance, in R if you want to create a website for your R package, running:

```
usethis::use_pkgdown()
usethis::use_pkgdown_travis()
```

...will set up the necessary infrastructure to create a website from your package documentation and to get Travis to automatically build, create, and deploy your package website to be put online. So for R projects, the `.travis.yml` file looks something like this after using the above functions:

```
language: R
sudo: false
cache: packages

before_cache: Rscript -e 'remotes::install_cran("pkgdown")'
deploy:
  provider: script
  script: Rscript -e 'pkgdown::deploy_site_github()'
  skip_cleanup: true
```

The `deploy` key is the important one here, as it is the one that tells Travis to send the generated website files to be put online (via GitHub Pages).

TODO: Add information for Python projects



## 23.7 Summary

FIXME: create concept map for integration.

## 23.8 Exercises

TODO: Get feedback on this and fill it out more later.

1. Setup Travis CI for the package you've been working on (for the course).
2. Get Travis to run unit tests of your package (not applicable for R packages).
3. Write a new unit test in your package that you *know* will fail. Push the changes up to GitHub.
4. Go through the Travis CI log and see what the log says about the failure and how to fix it. Since you made it purposefully fail, you will already know how to fix it, but the point is to get comfortable looking through Travis' logs.
5. Correct the test, push to GitHub, and get Travis to build your package successfully.

## 23.9 Key Points



## Chapter 24

# Publishing

### 24.1 Questions

### 24.2 Objectives

### 24.3 Introduction

This lesson looks at what should be included in and/or alongside reports and how best to do that. We use the generic term “report” to include research papers, summaries for clients, and everything else that is shorter than a book and is going to be read by someone else.

Our motivation is summed up in this quotation:

An article about computational science in a scientific publication is *not* the scholarship itself, it is merely *advertising* of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures.

– Jonathan Buckheit and David Donoho, paraphrasing Jon Claerbout, in Buckheit and Donoho (1995)

As the quote suggests, modern publishing involves much more than simply producing a report. It involves providing readers with the data underpinning the report, as well as any code written in analysing the data.

While the definition of data (and its associated metadata) is relatively easy to understand, code can come in many different forms. Here we distinguish between “analysis scripts” written solely for the purpose of the report (e.g. to

produce a figure) and “analysis software” that is formally packaged and released for use by a wider audience. Of course, in reality the analysis scripts/software written during the process of preparing a report can often lie on a continuum between those two definitions.

While some reports, datasets, software packages and/or analysis scripts can’t be published without violating personal or commercial confidentiality, every researcher’s default should be to make all these components of their work as widely available as possible. That means publishing it in an open access venue (Chapter 18) so that people who aren’t in academia can find and access it.

## 24.4 Identification

Before publishing anything, we need to understand the systems used to identify works and their authors.

A Digital Object Identifier (DOI) is a unique identifier for a particular version of a particular digital artifact such as a report, a dataset, or a piece of software. DOIs are written as `doi:prefix/suffix`, but you will often also see them represented as URLs like `http://dx.doi.org/prefix/suffix`. In order to be allowed to issue a DOI, online platforms (e.g. academic journals, data archives) must guarantee a certain level of security, longevity and access.

An ORCID is an Open Researcher and Contributor ID. You can get an ORCID for free, and you should include it in publications because people’s names and affiliations change over time.

## 24.5 Publishing a report

The best option for publishing a report (with a platform that issues a DOI) depends on the context. For academic, peer-reviewed research papers, numerous open access journals have popped up in recent years. Many formerly closed-access journals also now offer an open access option (for an additional fee). Another option is to publish with an online pre-print server (e.g. bioRxiv, arXiv). A preprint is a version of an academic research paper that precedes formal peer review and publication in a peer-reviewed journal. The preprint may be available, often as a non-typeset version available free, before and/or after a paper is published in a journal.

Online writing platforms such as Authorea are also an option, which allow a report to be openly viewable on the web throughout the entire writing process. At the end of the process Authorea can issue its own DOI for the report, or the text can be exported in the format required for submission to an academic journal. Finally, online platforms such as Figshare and Zenodo are a place

where any research outputs (reports, datasets, code, supplementary figures) can be published with a DOI. It is common for people to upload reports to these platforms so that they can be easily accessed by others.

## 24.6 Publishing data

The first step in publishing the data associated with a report is to determine what (if anything) needs to be published.

If the report involved the analysis of a publicly available dataset that is maintained and documented by a third party (e.g. open government data), then it's likely that no data publishing is required. The report simply needs to document where to access the data and what version was analyzed, along with any scripts and software used to download and process the data. In other words, it's not necessary to re-publish a duplicate of the original dataset if it's already accessible elsewhere.

Strictly speaking, it's not necessary to publish any data files produced during the analysis of a publicly available dataset either, since readers have access to the original data and the scripts/software used to process it. Having said that, it can be advantageous to publish processed data that is difficult to reproduce (e.g. it might require access to a supercomputing facility to run the code) and/or represents a derived quantity with high re-use potential. For instance, it would be worthwhile to publish an estimate of the global average surface temperature derived from a large database of weather observations, because a simple metric of global warming could be useful in many subsequent studies.

If a report involves the generation of a new dataset (e.g. observations collected during a field experiment), then clearly that dataset needs to be published. This section describes how to go about doing that.

### 24.6.1 What is the most useful way to share my data?

Making data useful to other people (including your future self) is one of the best investments you can make. The simple version of how to do this is:

- Always use tidy data.
- Include keywords describing the data in the project's `README.md` so that they appear on its home page and can easily be found by search engines.
- Give every dataset and every report a unique identifier (Section 24.4).
- Put data in open repositories.
- Use well-known formats like CSV and HDF5.
- Include an explicit license in every project and every dataset.
- Include units and other metadata.

The last point is often the hardest for people to implement, since many researchers have never seen a well-documented dataset. We draw inspiration from the data catalog included in the repository for the article “Women’s Pockets Are Inferior” and include a file `./data/README.md` in every project that looks like this:

- Infants born to women with HIV receiving an HIV test within two months of birth, 2017
  - ``Infant_HIV_Testing_2017.xlsx``
    - What is this?: Excel spreadsheet with summarized data.
    - Source(s): UNICEF, <<https://data.unicef.org/resources/dataset/hiv-aids-status-report-2017>>
    - Last Modified: July 2018 (according to website)
    - Contact: Greg Wilson <[greg.wilson@rstudio.com](mailto:greg.wilson@rstudio.com)>
    - Spatial Applicability: global
    - Temporal Applicability: 2009-2017
  - ``infant_hiv.csv``
    - What is this?: CSV export from ``Infant_HIV_Testing_2017.xlsx``
  - Notes
    - Data is not tidy: some rows are descriptive comments, others are blank separator rows
    - Use ``tidy_infant_hiv()`` to tidy this data.
- Maternal health indicators disaggregated by age
  - ``maternal_health_adolescents_indicators_April-2016_250d599.xlsx``
    - What is this?: Excel spreadsheet with summarized data.
    - Source(s): UNICEF, <<https://data.unicef.org/resources/dataset/maternal-health-report-2016>>
    - Last Modified: July 2018 (according to website)
    - Contact: Greg Wilson <[greg.wilson@rstudio.com](mailto:greg.wilson@rstudio.com)>
    - Spatial Applicability: global
    - Temporal Applicability: 2000-2014
  - ``at_health_facilities.csv``
    - What is this?: percentage of births at health facilities by country, year, and mode of delivery
    - Source(s): single sheet from ``maternal_health_adolescents_indicators_April-2016_250d599.xlsx``
  - ``c_sections.csv``
    - What is this?: percentage of Caesarean sections by country, year, and mode of delivery
    - Source(s): single sheet from ``maternal_health_adolescents_indicators_April-2016_250d599.xlsx``
  - ``skilled_attendant_at_birth.csv``
    - What is this?: percentage of births with skilled attendant present by country, year, and mode of delivery
    - Source(s): single sheet from ``maternal_health_adolescents_indicators_April-2016_250d599.xlsx``
  - Notes
    - Data is not tidy: some rows are descriptive comments, others are blank separator rows
    - Use ``tidy_maternal_health_adolescents()`` to tidy this data.

The catalog above doesn’t include column headers or units because the data isn’t tidy. It *does* include the names of the functions used to reformat that data, and `./results/README.md` then includes the information that users will want. One section of that file is shown below:

- Infants born to women with HIV receiving an HIV test within two months of birth, 2017
  - `infant_hiv.csv`

- What is this?: tidied version of CSV export from spreadsheet.
- Source(s): UNICEF, <<https://data.unicef.org/resources/dataset/hiv-aids-statistical-table>>
- Last Modified: September 2018
- Contact: Greg Wilson <[greg.wilson@rstudio.com](mailto:greg.wilson@rstudio.com)>
- Spatial Applicability: global
- Temporal Applicability: 2009-2017
- Generated By: scripts/tidy-24.R

| Header   | Datatype | NA    | Description                                 |
|----------|----------|-------|---|
| country  | char     | false | ISO3 country code of country reporting data |
| year     | integer  | false | year CE for which data reported             |
| estimate | double   | true  | estimated percentage of measurement         |
| hi       | double   | true  | high end of range                           |
| lo       | double   | true  | low end of range                            |

Note that this catalog includes both units and whether or not a field can be NA. Note also that calling a field “NA” is asking for trouble...

## 24.6.2 What standards of data sharing should I aspire to?

The FAIR Principles describe what research data should look like. They are still aspirational for most researchers, but they tell us what to aim for. The most immediately important elements of the FAIR Principles are outlined below.

### 24.6.2.1 Data should be *findable*.

The first step in using or re-using data is to find it. You can tell you’ve done this if:

1. (Meta)data is assigned a globally unique and persistent identifier (Section 24.4).
2. Data is described with rich metadata (like the catalog shown above).
3. Metadata clearly and explicitly includes the identifier of the data it describes.
4. (Meta)data is registered or indexed in a searchable resource, such as the data sharing platforms described in Section 24.6.

### 24.6.2.2 Data should be *accessible*.

You can’t use data if you don’t have access to it. In practice, this rule means the data should be openly accessible (the preferred solution) or that authenticating in order to view or download it should be free. You can tell you’ve done this if:

1. (Meta)data is retrievable by its identifier using a standard communications protocol like HTTP.
2. Metadata is accessible even when the data is no longer available.

#### 24.6.2.3 Data should be *interoperable*.

Data usually needs to be integrated with other data, which means that tools need to be able to process it. You can tell you've done this if:

1. (Meta)data uses a formal, accessible, shared, and broadly applicable language for knowledge representation
2. (Meta)data uses vocabularies that follow FAIR principles
3. (Meta)data includes qualified references to other (meta)data

#### 24.6.2.4 Data should be *reusable*.

This is the ultimate purpose of the FAIR Principles and much other work. You can tell you've done this if:

1. Meta(data) is described with accurate and relevant attributes.
2. (Meta)data is released with a clear and accessible data usage license.
3. (Meta)data has detailed provenance.
4. (Meta)data meets domain-relevant community standards.

### 24.6.3 How and where do I publish the data?

Small datasets (i.e., anything under 500 MB) can be stored in version control using the conventions described in Chapter 17. If the data is being used in several projects, it may make sense to create one repository to hold only the data; the R community refers to these as data packages, and they are often accompanied by small scripts to clean up and query the data. Be sure to give the dataset an identifier as discussed in Section 24.4.

For medium-sized datasets (between 500 MB and 5 GB), it's better to put the data on platforms like the Open Science Framework, Dryad, and Figshare. Each of these will give the datasets identifiers; those identifiers should be included in reports along with scripts to download the data. Big datasets (i.e., anything more than 5 GB) may not be yours in the first place, and probably need the attention of a professional archivist. Any processed or intermediate data that takes a long time to regenerate should probably be published as well using these same sizing rules; all of this data should be given identifiers, and those identifiers should be included in reports.

#### Data journals



While archiving data at a site like Dryad or Figshare (following the FAIR Principles) is usually the end of the data publishing process, there is the option of publishing a journal paper to describe the dataset in detail. Some research disciplines have journals devoted to describing particular types of data (e.g. Geoscience Data Journal) and there are also generic data journals (e.g. Scientific Data).

## 24.7 Publishing analysis software

In the preceding chapters we have learned how to document and package software so that it can be installed and used by others. The final step in this process is publication.

It is common practice to have the code associated with a software package openly available on a hosting service such as GitHub (or GitLab or Bitbucket). These hosting services are not only a convenient place for people to ask questions and make contributions/improvements to the software, they also have built-in functionality for managing the release of new versions of the software. One limitation of these sites, however, is that they don't guarantee persistent long term storage (e.g. if you changed the name of your GitHub repository any URLs for the existing repository would be broken). Acknowledging this limitation, GitHub provides Zenodo integration for creating a DOI with each new software release.

### Software journals

While creating a DOI using a site like Zenodo is often the end of the software publishing process, there is the option of publishing a journal paper to describe the software in detail. Some research disciplines have journals devoted to describing particular types of software (e.g. Geoscientific Model Development), and there are also a number of generic software journals such as the Journal of Open Research Software and Journal of Open Source Software.

## 24.8 Publishing analysis scripts

The final component that needs to be published is the analysis scripts. Unlike analysis software that has been packaged and released for use by a wider audience, analysis scripts are simply written to create the figures and tables presented in a given report. In fact, these scripts would typically make use of analysis software written by the wider data science community (e.g. matplotlib, ggplot) as well discipline-specific packages written by colleagues or co-authors (e.g. AstroPy).

Given that analysis scripts will typically leverage a wide variety of existing software packages, there's actually three separate items that need to be published:

1. A detailed description of the analysis software used
2. A copy of any analysis scripts written by the authors to produce the key results presented in the report
3. A description of the data processing steps taken in producing each key result (i.e. a step-by-step account of how the software and scripts were actually implemented)

Earlier we saw that there are well-developed and widely adopted guidelines for data publishing (e.g. the FAIR Principles). The same is not true for analysis scripts. Librarians, publishers, and regulatory bodies are still trying to determine the best way for code to be documented and archived. For the moment, the best advice we can give for those three key items is discussed below. That advice ranges from the bare minimum that needs to be done through to current gold standard practice.

### 24.8.1 Software description

In order to document the software packages that were used, the bare minimum requirement is to list the name and version number of each software package that played a critical role in producing the analysis presented in your report.

As Section 22.6 described, you can get these automatically by running:

```
$ pip freeze > requirements.txt
```

For everything else, you should write a script or create a rule in your project's Makefile (Chapter 14), since the commands used to get version numbers will vary from tool to tool:

```
## versions : dump versions of software.
versions :
    @echo '# Python packages'
    @pip freeze
    @echo '# dezply'
    @dezply --version
    @echo '# parajune'
    @parajune --status | head 1
```

While such a list means your software environment is now technically reproducible, you've left it up to the reader to figure out how to get all those software packages and libraries installed and playing together nicely. In some cases this is fine (e.g. it might be easy enough for a reader to install the handful of R packages you used), but in other cases you might want to save the reader (and your future self) the pain of software installation by making use of a tool that can automatically install a specified software environment. The most prominent

such tool in data science at the moment is conda. A conda environment can be exported,

```
$ conda env export -n myenv -f myenv.yml
```

and made available so that readers can use it to install the same environment on their own computer:

```
$ conda env create -f myenv.yml
```

### Conda environments

The Python for Atmosphere and Ocean Scientists lesson materials maintained by Data Carpentry have a section devoted to Software Installation using Conda

Beyond conda there are more complex tools like Docker and Nix, which can literally install your entire environment (down to the precise operating system) on a different computer. There's lots of debate about the potential and suitability of these tools as a solution to reproducible research, but it's fair to say that their complexity puts them out of reach for many researchers.

## 24.8.2 Analysis scripts

The next item you'll need to publish is a copy of the scripts written to execute those software packages. Depending on the size or complexity of the scripts you have written, and whether you re-use them in multiple projects, you may publish script by script or create a zip file or tar file that includes everything. For example, the Makefile fragment below creates `~/archive/meow-2019-02-21.tgz`:

```
ARCHIVE=${HOME}/archive
PROJECT=meow
TODAY=$(shell date "+%Y-%m-%d")
SCRIPTS=./Makefile ./bin/*.py ./bin/*.sh

## archive : create an archive of all the scripts used in this run
archive :
    @mkdir -p ${ARCHIVE}
    @tar zcf ${ARCHIVE}/${PROJECT}-${TODAY}.tgz
```

## 24.8.3 Data processing steps

A software description and analysis scripts on their own are not much use to a reader; they also need to know how those scripts was actually executed. This means including the configuration files (Chapter 13), and/or command-line parameters used to generate each key result.

The way in which this information is collected and archived depends on how your workflow is constructed. If all of a program’s parameters are in a configuration file (Chapter 13), then that file can be archived. Alternatively, you might need to have your program print out its configuration parameters and then use `grep` or a script to extract information from the logfile (Chapter E.3.5).

If your workflow involves executing a series of command line programs, then you can keep a log/record of the command line entries required to produce a given result. For example, the `cmdline-provenance` package generates such records, including keeping track of the corresponding version control revision number, so you know exactly which version of your command line program was executed.

As before, while these bare minimum log files ensure that your workflow is reproducible, they may not be particularly comprehensible. Manually recreating workflows from them might be a tedious and time consuming process, even for just moderately complex analyses. To make things a little easier for the reader (and your future self), it’s a good idea to include a README file in your code library explaining the sequence of commands required to produce common/key results. You might also provide a Makefile that automatically builds and executes common workflows. Beyond that the options get more complex, with workflow management packages like VisTrails providing a graphical interface that allows users to drag and drop the various components of their workflow.

#### 24.8.4 Where to publish all this stuff?

Following the steps above, you’ll be left with a text file (or perhaps an environment file exported from a conda) describing your software environment, a copy of your code library and various log files, and README files and/or Makefiles that describe your data processing steps. Sites like Figshare and Zenodo are the ideal place to publish these items, as they have been specifically setup for archiving the “long tail” of reports (e.g. supplementary figures, tables, code and data).

### 24.9 Summary

FIXME: create concept map for publishing

#### 24.10 Exercises

##### ORCID

If you don’t already have an ORCID, go to the website and register now.

If you do have an ORCID, login at the website and make sure that your details and publication record are up-to-date.

**A FAIR test**

An online questionnaire for measuring the extent to which datasets are FAIR has been created by the Australian Research Data Commons.

Take the questionnaire for a dataset you have published or that you use often.

**Publishing your code**

Think about a project that you're currently working on.

How would you go about publishing the code associated with that project? (i.e. the software description, analysis scripts and data processing steps)

## 24.11 Key Points



## Chapter 25

# Finale

### 25.1 Questions

- What have we learned?

### 25.2 Objectives

### 25.3 Introduction

Ever since Wing (2006) introduced the term “computational thinking” in 2006, computer scientists and educators have debated what exactly it means Denning (2017). What isn’t debated is the fact that programmers tend to think about problems in ways inspired by programming. This paper describes ten of those ways which may help people in other domains see their own problems with fresh eyes.

*An early version of this paper was inspired by Jon Udell’s essay “Seven Ways to Think Like the Web” Udell (2011).*

### 25.4 Programs are data.

The key insight that all of modern computing is built on is that programs are just another kind of data. Source code is text, no different from a thesis or a poem; it can be searched with the same tools used to search other kinds of documents, and new text (such as documentation) can be generated from old.

More importantly, once a program is loaded into memory its instructions are just bytes, no different in principle or practice from those used to represent the pixels in an image or the numbers in a vector field. Those blocks of instructions can be stored in data structures, passed to functions, or altered on the fly, just like other data. Almost all advanced programming techniques depend on this fact, from function pointers in C and callbacks in JavaScript to decorators in Python and lazy evaluation in R.

## 25.5 Computers don't understand anything.

Computers don't understand: they only obey instructions that make them appear to. If a person looks at Figure ??, they see the word “data”:

A machine doesn't; it doesn't even see four blobs of blue pixels on a gray background, because it doesn't “see” anything. Equally, calling a variable “temperature” doesn't mean the computer will store a temperature in it—it would do exactly the same thing if the variable was called “pressure” or “frankenstein” or “a7”. This may seem obvious once you have written a few programs, but thirty years after Pea (1986) first called it the “superbug”, believing that the computer will somehow understand intent remains a common error.

## 25.6 Programming is about creating and combining abstractions.

Computers don't have to understand instructions in order to execute them, but we do in order to create them (and fix them afterward). Since our working memory can only juggle a handful of things at once Miller (1956), we have to create abstractions so that our representations of problems will fit into hardware whose performance doubles over millions of years rather than every eighteen months.

The key to making workable abstractions is to separate “what” from “how”, or in computing terms, to separate *interface* and *implementation*. An interface is what something knows how to do; its implementation is how it does those things. There can be dozens of different implementations of a single interface: if we do our work well, we shouldn't have to care about the differences between them until something goes wrong or we need to improve its performance.



## 25.7 Every redundancy in software is an abstraction trying to be born.

The history of programming is in part the history of people noticing patterns and then making it easier for programmers to use them. Does your program repeatedly search an array to find the largest and smallest values? Write a function called `bounds`. Does it repeatedly search arbitrary data structures to find values that meet certain criteria? Write a generator that returns values one by one and another function that filters those according to some criteria.

The problem with eliminating redundancy is that it can make software denser, which in turn makes it harder for non-specialists to understand. Like mathematics and modern art, it can take years of training for someone to reach the point where they can see how beautiful something is.

The other problem with patterns is that they can lead to expert blind spot Nathan and Petrosino (2003). Once experts have internalized patterns, they are often unable to remember that they ever saw the world any other way, or to see the world afresh through novice eyes. As a result, they are prone to say, “It’s obvious,” and then follow it with incomprehensible jargon.

## 25.8 Create models for computers and views for human beings.

One consequence of the second and third rules is important enough to be a rule in its own right: we should create models for computers and views for human beings. A *model* is a precise, detailed representation of information that is easy for a computer to operate on; a *view* is a way of displaying information that human beings can easily understand and interact with. For example, an HTML page is represented in memory as a data structure containing nodes for elements like headings and paragraphs, which can in turn contain a mix of other nodes or text (Figure ??).

That model can be rendered in a browser, turned into speech, or displayed as text using angle brackets. None of these *is* the model: they’re all views that make the information in the model comprehensible in different ways.

Turning a model into a view is hard, but turning a view back into a model is harder. For example, parsing the textual representation of HTML takes thousands of lines of code, but doing OCR or speech recognition to translate a rendered page or its spoken equivalent back into nodes and text can take millions. Ironically, the same programmers who insist on this separation in software they build for the rest of humanity have been remarkably resistant to the idea of adopting any kind of model-view separation in programming itself.

## 25.9 Paranoia makes us productive.

“I want to count the cells in this photograph” is easy to say, but what does it actually mean? What constitutes a cell? When do you decide that a lumpy blob of pixels is two cells rather than one, or three instead of two? Every program embodies decisions like these; the sooner these decisions are made explicit and the earlier they’re checked, the more productive we will be. The precise order of steps doesn’t seem to matter: we can write tests, then write software to make those tests pass, or write the program first and then test it. What *does* matter is alternating development and testing in short interleaved bursts Fucci et al. (2017).

Of course, we don’t stop worrying once we’ve typed our code in. We check that data is formatted properly to protect ourselves against “garbage in, garbage out”. We put checks in our code to make sure that parameters are sensible, data structures consistent, files aren’t empty, and so on. This is called *defensive programming*, and one of the signs of a mature programmer is a high density of assertions and other self-checks in her code. It is harder to do this in research software than in most commercial software because (almost by definition) we don’t know what the right answer *is* in research software, which makes it difficult to check that we’re getting it.

## 25.10 Things that don’t change are easier to understand than things that do.

Programmers use the words *mutable* and *immutable* to refer to things that can be modified after they are created and things that cannot. Immutable things are easier to understand because you don’t have to re-trace their history in order to understand their state. However, immutable data can be less efficient than mutable data: for example, it’s very expensive to make a copy of an entire multi-megabyte image just because we want to change one pixel.

Older programming languages like C and Fortran allowed most data to be mutable because computer time was expensive. This led to many hard-to-find bugs, so newer languages either make data immutable or automatically copy data when asked to make changes in order to give the appearance of immutability.

One special case of this rule is automating workflows. As Whitehead (1958) said, “Civilization advances by extending the number of important operations which we can perform without thinking about them.” Every time we automate a task—i.e., make its steps immutable—we reduce the chances of getting it wrong the next time, and have more time to think about things that machines *can’t* do for us.

## 25.11 Better algorithms are better than better hardware.

One of the greatest mathematical advances of the Twentieth Century was the idea of *algorithmic complexity*. The key insight is that the running time or memory requirements of an algorithm grow in a predictable way as a function of the size of the problem we are trying to solve Conery (2016). Some algorithms slow down gently as their inputs get larger, but others slow down so much that even if the whole universe was one large computer, it couldn't solve any interesting problem. Faster chips are therefore very welcome, but the real key to speed is to focus on how we're doing things.

## 25.12 Distributed is different.

Distributed computing is intrinsically different from running a program on a single machine Waldo et al. (1994). On a single computer, we can usually pretend that nobody else is modifying the data while we're trying to use it. Once our data is distributed, that simplification breaks down, and we have to worry about things like someone else adding records to the database between the time we ask how many there are and the time we start processing them.

Similarly, we can pretend that a program running on a single computer either works or doesn't. When that same program is accessing remote resources, we have to worry about whether a long delay means that something has failed, or whether it's just being slow. Many attempts have been made to paper over these differences, but all have failed in the large. As a result, the future of programming is about how we deal with this—a statement that has been true since the 1980s

## 25.13 Privacy, security, fairness, and responsibility can't be added after the fact.

Our final rule may be the most important of all. As the last few years have shown, collecting and interpreting data is never a neutral activity: who we share data with, how we classify it, and most importantly, who gets to decide these things are all political decisions with ever-increasing impact, and we are past the point where we can pretend otherwise.

Attempts to add privacy, security, and fairness to systems after they have been built and deployed have repeatedly failed. The other “future of programming” is therefore to take digital health as seriously as physical health, and to make

those responsible for it as accountable as those responsible for other aspects of our wellbeing.

## 25.14 Summary

Artisans have known for centuries that the tool shapes the hand. If computers are tools for thinking with, then it shouldn't surprise us that writing software shapes the minds of those who do it. The ten rules listed above are just a few reflections of this; we hope that they, and the skills and tools taught in these lessons, will help you see the world in new ways.

## 25.15 Exercises

FIXME: exercise for finale

## 25.16 Key Points

# Appendix A

## License

*This is a human-readable summary of (and not a substitute for) the license. Please see <https://creativecommons.org/licenses/by/4.0/legalcode> for the full legal text.*

This work is licensed under the Creative Commons Attribution 4.0 International license (CC-BY-4.0).

### You are free to:

- **Share**—copy and redistribute the material in any medium or format
- **Remix**—remix, transform, and build upon the material for any purpose, even commercially.

The licensor cannot revoke these freedoms as long as you follow the license terms.

### Under the following terms:

- **Attribution**—You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
- **No additional restrictions**—You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

### Notices:

You do not have to comply with the license for elements of the material in the public domain or where your use is permitted by an applicable exception or limitation.

No warranties are given. The license may not give you all of the permissions necessary for your intended use. For example, other rights such as publicity, privacy, or moral rights may limit how you use the material.

## Appendix B

# Code of Conduct

In the interest of fostering an open and welcoming environment, we as contributors and maintainers pledge to making participation in our project and our community a harassment-free experience for everyone, regardless of age, body size, disability, ethnicity, gender identity and expression, level of experience, education, socio-economic status, nationality, personal appearance, race, religion, or sexual identity and orientation.

### B.1 Our Standards

Examples of behavior that contributes to creating a positive environment include:

- using welcoming and inclusive language,
- being respectful of differing viewpoints and experiences,
- gracefully accepting constructive criticism,
- focusing on what is best for the community, and
- showing empathy towards other community members.

Examples of unacceptable behavior by participants include:

- the use of sexualized language or imagery and unwelcome sexual attention or advances,
- trolling, insulting/derogatory comments, and personal or political attacks,
- public or private harassment,
- publishing others' private information, such as a physical or electronic address, without explicit permission, and
- other conduct which could reasonably be considered inappropriate in a professional setting

## B.2 Our Responsibilities

Project maintainers are responsible for clarifying the standards of acceptable behavior and are expected to take appropriate and fair corrective action in response to any instances of unacceptable behavior.

Project maintainers have the right and responsibility to remove, edit, or reject comments, commits, code, wiki edits, issues, and other contributions that are not aligned to this Code of Conduct, or to ban temporarily or permanently any contributor for other behaviors that they deem inappropriate, threatening, offensive, or harmful.

## B.3 Scope

This Code of Conduct applies both within project spaces and in public spaces when an individual is representing the project or its community. Examples of representing a project or community include using an official project e-mail address, posting via an official social media account, or acting as an appointed representative at an online or offline event. Representation of a project may be further defined and clarified by project maintainers.

## B.4 Enforcement

Instances of abusive, harassing, or otherwise unacceptable behavior may be reported by emailing the project team. All complaints will be reviewed and investigated and will result in a response that is deemed necessary and appropriate to the circumstances. The project team is obligated to maintain confidentiality with regard to the reporter of an incident. Further details of specific enforcement policies may be posted separately.

Project maintainers who do not follow or enforce the Code of Conduct in good faith may face temporary or permanent repercussions as determined by other members of the project's leadership.

## B.5 Attribution

This Code of Conduct is adapted from the Contributor Covenant version 1.4.



## Appendix C

# Contributing

Contributions of all kinds are welcome, from errata and minor improvements to entirely new sections and chapters: please email us or submit an issue or pull request to our GitHub repository. Everyone whose work is incorporated will be acknowledged; please note that all contributors are required to abide by our Code of Conduct. See the contributing guidelines for more detail.

Please note that we use Simplified English rather than Traditional English, i.e., American rather than British spelling and grammar. We encourage translations; if you would like to take this on, please email us.

If you wish to report errata or suggest improvements to wording, please include the chapter name in the first line of the body of your report (e.g., **Testing Data Analysis**).



# Appendix D

## Glossary

**Absolute error** FIXME

**Accuracy** FIXME

**Action (in Make)** FIXME

**Active listening** FIXME

**Actual output (of a test)** FIXME

**Agile development** FIXME

**Ally** FIXME

**Analysis and estimation** FIXME

**Annotated tag (in version control)** FIXME

**Append mode** FIXME

**Application Programming Interface (API)** FIXME

**Assertion** FIXME

**Authentic task** A task which contains important elements of things that learners would do in real (non-classroom situations). To be authentic, a task should require learners to construct their own answers rather than choose between provided answers, and to work with the same tools and data they would use in real life.

**Auto-completion** FIXME

**Automatic variable** FIXME

**Backlog** FIXME

**Binary code** FIXME

**Bit rot** FIXME

**Branch-per-feature workflow** branch-per-feature

**Breakpoint** FIXME

**Buffer** FIXME

**Build tool** FIXME

**Byte code** FIXME

**Call stack** FIXME

**Camel case** FIXME

- Catch (an exception)** FIXME
- Code browser** FIXME
- Cognitive load** FIXME
- Comma-separated values (CSV)** FIXME
- Commit hash** FIXME
- Commit message** FIXME
- Commons** FIXME
- Competent practitioner** Someone who can do normal tasks with normal effort under normal circumstances. See also novice and expert.
- Computational competence** FIXME
- Computational stylometry** FIXME
- Computational thinking** FIXME
- Conditional expression** FIXME
- Configuration object** FIXME
- Context manager** FIXME
- Continuous integration** FIXME
- Coverage** FIXME
- Creative Commons - Attribution License (CC-BY)** FIXME
- Data engineering** FIXME
- Data package** FIXME
- Declarative programming** FIXME
- Default target** FIXME
- Dependency graph** FIXME
- Design patterns** FIXME
- Destructuring** FIXME
- Dictionary** FIXME
- Digital Object Identifier (DOI)** FIXME
- Directory** A folder in a filesystem.
- Directory Tree** If the nesting relationships between directories in a filesystem are drawn as arrows from the containing directory to the nested ones, a tree structure develops.
- Disk** Disk refers to disk storage, a physical component of a computer that stores information on a disk. The most common kind of disk storage is a hard disk drive, which is a storage drive with a non-removable disk.
- Docstring** FIXME
- Environment** FIXME
- Embedded documentation** FIXME
- Error (result from a unit test)** FIXME
- Escape sequence** FIXME
- Exception** FIXME
- Exit status** FIXME
- Expected output (of a test)** FIXME
- Expert** Someone who can diagnose and handle unusual situations, knows when the usual rules do not apply, and tends to recognize solutions rather than reasoning to them. See also competent practitioner and novice.
- Exploratory programming** FIXME

**Exponent** FIXME  
**External error** FIXME  
**Failure (result from a unit test)** FIXME  
**False beginner** Someone who has studied a language before but is learning it again. False beginners start at the same point as true beginners (i.e., a pre-test will show the same proficiency) but can move much more quickly.  
**Feature boxing** FIXME  
**Feature creep** FIXME  
**Filename stem** FIXME  
**Filesystem** Controls how files are stored and retrieved on disk by an operating system. Also used to refer to the disk that is used to store the files or the type of the filesystem.  
**Fixture** FIXME  
**Flag** FIXME  
**Forge** FIXME  
**Fork** FIXME  
**Format string** FIXME  
**Frequently Asked Questions (FAQ)** FIXME  
**Fully-qualified name** FIXME  
**Function (in Make)** FIXME  
**Function attribute** FIXME  
**GitHub Pages** FIXME  
**Globbering** FIXME  
**GNU Public License (GPL)** FIXME  
**Hot spot** FIXME  
**HTTP status code** FIXME  
**Impostor syndrome** FIXME  
**In-place operator** FIXME  
**Index** FIXME  
**Install** FIXME  
**Instrumenting profiler** FIXME  
**Integrated Development Environment (IDE)** FIXME  
**Internal error** FIXME  
**ISO date format** FIXME  
**Issue** FIXME  
**Iteration (in software development)** FIXME  
**Jenny (a repository)** FIXME  
**Join (of database tables)** FIXME  
**JSON** FIXME  
**Kebab case** FIXME  
**Label (in issue tracker)** FIXME  
**Library** FIXME  
**Lint** FIXME  
**List comprehension** FIXME  
**Logging framework** FIXME  
**Macro** FIXME

**Magnitude** FIXME

**Makefile** FIXME

**Mantissa** FIXME

**Memory** A physical device on your computer that temporarily stores information for immediate use.

**Mental model** A simplified representation of the key elements and relationships of some problem domain that is good enough to support problem solving.

**Method** A function that is specific to an object type, based on qualities of that type, e.g. a string method like `upper()` which turns characters in a string to uppercase.

**MIT License** FIXME

**Not Invented Here (NIH)** FIXME

**Novice** Someone who has not yet built a usable mental model of a domain. See also competent practitioner and expert.

**Object** An object is a programming language's way of describing and storing values, usually labeled with a variable name.

**Object-oriented programming** FIXME

**Open science** FIXME

**Operational test** FIXME

**Oppression** FIXME

**ORCID** FIXME

**Overlay configuration** FIXME

**Package** FIXME

**Pair programming** FIXME

**Pattern rule** FIXME

**Peer action** FIXME

**Phony target** FIXME

**Post-mortem** FIXME

**Pothole case** FIXME

**Precision** FIXME

**Prerequisite (in Make)** FIXME

**Privilege** FIXME

**Procedural programming** FIXME

**Product manager** FIXME

**Profiler** FIXME

**Project manager** FIXME

**Provenance** FIXME

**Pseudorandom number generator (PRNG)** FIXME

**Public domain license (CC-0)** FIXME

**Pull request** FIXME

**Raise** FIXME

**Raster image** FIXME

**Rebase** FIXME

**Redirection** FIXME

**Refactor** FIXME

Relative error FIXME  
Remote login FIXME  
Repository FIXME  
Representation State Transfer (**REST**) FIXME  
Reproducible example (**reprx**) FIXME  
Reproducible research FIXME  
Research software engineer (**RSE**) FIXME  
reStructured Text FIXME  
Rotating file FIXME  
Rule (in Make) FIXME  
Sampling profiler FIXME  
Scalable Vector Graphics (**SVG**) FIXME  
Seed (for pseudorandom number generator) FIXME  
Semantic versioning FIXME <https://semver.org/>  
Set and override (**pattern**) FIXME  
Shebang FIXME  
Short circuit test FIXME  
Side effects FIXME  
Sign FIXME  
Silent error FIXME  
Silent failure FIXME  
Situational action FIXME  
Software development process FIXME  
Source code FIXME  
SSH key FIXME  
SSH protocol FIXME  
Stand-up meeting FIXME  
Standard error FIXME  
Standard input FIXME  
Standard output FIXME  
Streaming data FIXME  
Sturdy development FIXME  
Subsampling FIXME  
Success (result from a unit test) FIXME  
Sustainability FIXME  
Symbolic debugger FIXME  
Syntax highlighting FIXME  
Synthetic data FIXME  
Tab completion FIXME  
Tag (in version control) FIXME  
Target (in Make) FIXME  
Target (of oppression) FIXME  
Technical debt FIXME  
Test coverage FIXME  
Test-driven development FIXME  
Test framework FIXME

**Test runner** FIXME

**Tidy data** As defined in Wick2014, tabular data is tidy if: 1. Each variable is in one column. 2. Each different observation of that variable is in a different row. 3. There is one table for each kind of variable. 4. If there are multiple tables, each includes a key so that related data can be linked.

**Time boxing** FIXME

**Timestamp (on a file)** FIXME

**Tolerance** FIXME

**Triage** FIXME

**Tuning** FIXME

**Tuple** FIXME

**Typesetting language** FIXME

**Unit test** FIXME

**Update operator** See in-place operator.

**Validation** FIXME

**Variable (in Make)** FIXME

**Variable (in Python)** A symbolic name that reserves memory to store a value.

**Vector image** FIXME

**Verification** FIXME

**Virtual environment** FIXME

**What You See Is What You Get (WYSIWYG)** FIXME

**Wildcard** FIXME

**Working directory** FIXME

**Working memory** FIXME

**Wrapper** FIXME

**YAML** FIXME



## Appendix E

# Learning Objectives

### E.1 Novice R

FIXME: fill in learning objectives for novice R

### E.2 Novice Python

FIXME: fill in learning objectives for novice Python



## **E.3 RSE**

**E.3.1 The Unix Shell**

**E.3.2 Automating Analyses**

**E.3.3 Syndicating Data**

**E.3.4 Configuring Software**

**E.3.5 Logging**

**E.3.6 Unit Testing**

**E.3.7 Verification**

**E.3.8 A Branching Workflow**

**E.3.9 Managing Backlog**

**E.3.10 Code Style and Review**

**E.3.11 Development Process**

**E.3.12 Continuous Integration**

**E.3.13 Working Remotely**

**E.3.14 Other Tools**

**E.3.15 Documentation**

**E.3.16 Refactoring**

**E.3.17 Project Structure**

**E.3.18 Including Everyone**

**E.3.19 Python Packaging**

**E.3.20 Publishing**

**E.3.21 Teamwork**

**E.3.22 Pacing**

**E.3.23 Finale**



## Appendix F

# Key Points

### F.1 Novice R

FIXME: fill in keypoints for novice R

### F.2 Novice Python

FIXME: fill in keypoints for novice Python



## **F.3 RSE**

**F.3.1 The Unix Shell**

**F.3.2 Automating Analyses**

**F.3.3 Syndicating Data**

**F.3.4 Configuring Software**

**F.3.5 Logging**

**F.3.6 Unit Testing**

**F.3.7 Verification**

**F.3.8 A Branching Workflow**

**F.3.9 Managing Backlog**

**F.3.10 Code Style and Review**

**F.3.11 Development Process**

**F.3.12 Continuous Integration**

**F.3.13 Working Remotely**

**F.3.14 Other Tools**

**F.3.15 Documentation**

**F.3.16 Refactoring**

**F.3.17 Project Structure**

**F.3.18 Including Everyone**

**F.3.19 Python Packaging**

**F.3.20 Publishing**

**F.3.21 Teamwork**

**F.3.22 Pacing**

**F.3.23 Finale**





## Appendix G

### The Rules

1. Be kind: all else is details.
2. A week of hard work can sometimes save you an hour of thought.
3. While the hand makes the tool, the tool shapes the hand.
4. Nothing in software development makes sense except in light of human psychology.
5. If you need to write a parser, you've done something wrong.
6. Always seed your random number generator, and always record the seed.
7. Break any of these rules rather than doing something awkward.



Aurora, V. et al. 2018. How to respond to code of conduct reports. Frame Shift Consulting.

A practical step-by-step guide to handling code of conduct issues.

Bettenburg, N. et al. 2008. What makes a good bug report? *In* Proc. 16th ACM SIGSOFT international symposium on foundations of software engineering - (SIGSOFT'08/FSE'16). ACM Press.

Reports a survey of developers on what makes for a good bug report.

Bollier, D. 2014. Think like a commoner: A short introduction to the life of the commons. New Society Publishers.

A short introduction to a widely-used model of governance.

Brookfield, S.D., and S. Preskill. 2016. The discussion book: 50 great ways to get people talking. Jossey-Bass.

Describes fifty different ways to get groups talking productively.

Brown, M.J. 2007. Building powerful community organizations: A personal guide to creating groups that can solve problems and change the world. Long Haul Press.

A practical guide to creating effective organizations in and for communities.

Buckheit, J.B., and D.L. Donoho. 1995. WaveLab and reproducible research. *In* Wavelets and statistics. Springer New York. 55–81.

An early and influential discussion of reproducible research.

Conery, R. 2016. The imposter's handbook. Big Machine, Inc.  
FIXME

Denning, P.J. 2017. Remaining trouble spots with computational thinking. *Communications of the ACM*. 60:33–39. doi:10.1145/2998438.

FIXME

Fogel, K. 2005. Producing open source software: How to run a successful free software project. O'Reilly Media.

The definite guide to managing open source software development projects.

Freeman, J. 1972. The tyranny of structurelessness. *The Second Wave*. 2.

Points out that every organization has a power structure: the only question is whether it's accountable or not.

Fucci, D. et al. 2017. A dissection of the test-driven development process: Does it really matter to test-first or to test-last? *IEEE Transactions on Software Engineering*. 43. doi:10.1109/tse.2016.2616877.  
 FIXME

Fucci, D. et al. 2016. An external replication on the effects of test-driven development using a multi-site blind analysis approach. *In Proc. 10th ACM/IEEE international symposium on empirical software engineering and measurement (ESEM'16)*. ACM Press.  
 The latest in a long line to find that test-driven development (TDD) has little or no impact on development time or code quality.

Goldberg, D. 1991. What every computer scientist should know about floating-point arithmetic. *ACM Computing Surveys*. 23. doi:10.1145/103162.103163.  
 FIXME

Haddock, S., and C. Dunn. 2010. Practical computing for biologists. Sinauer Associates.  
 FIXME

Hekman, D.R. et al. 2017. Does diversity-valuing behavior result in diminished performance ratings for non-white and female leaders? *Academy of Management Journal*. 60:771–797. doi:10.5465/amj.2014.0538.  
 FIXME

Lindberg, V. 2008. Intellectual property and open source: A practical guide to protecting code. O'Reilly Media.  
 A thorough dive into intellectual property issues related to open source software

Miller, G.A. 1956. The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological Review*. 63:81–97. doi:10.1037/h0043158.  
 FIXME

Morin, A., and J.U.A.P. Sliz. 2012. A quick guide to software licensing for the scientist-programmer. *PLoS Computational Biology*. 8. doi:10.1371/journal.pcbi.1002598.  
 A short introduction to software licensing for non-specialists.

Nathan, M.J., and A. Petrosino. 2003. Expert blind spot among pre-service teachers. *American Educational Research Journal*. 40:905–928. doi:10.3102/00028312040004905.  
 FIXME

Noble, W.S. 2009. A quick guide to organizing computational biology projects. *PLoS Computational Biology*. 5. doi:10.1371/journal.pcbi.1000424.  
 How to organize a small to medium-sized bioinformatics project.

Patil, P. et al. 2016. A statistical definition for reproducibility and replicability. doi:10.1101/066803.

Pea, R.D. 1986. Language-independent conceptual “bugs” in novice programming. *Journal of Educational Computing Research*. 2:25–36. doi:10.2190/689t-1r2a-x4w4-29j2.

FIXME

Perez De Rosso, S., and D. Jackson. 2013. What’s wrong with git? *In Proc. 2013 ACM international symposium on new ideas, new paradigms, and reflections on programming and software (Onward!’13)*.

The first in a series exploring what’s wrong with Git.

Perez De Rosso, S., and D. Jackson. 2016. Purposes, concepts, misfits, and a re-design of git. *In Proc. 2016 ACM SIGPLAN international conference on object-oriented programming, systems, languages, and applications (OOPSLA’2016)*.

The second in a series exploring what’s wrong with Git.

Smith, P. 2011. *Software build systems: Principles and experience*. Addison-Wesley Professional.

A thorough, readable exploration of how software build systems and tools work.

Steinmacher, I. et al. 2014. The hard life of open source software project newcomers. *In Proc. 7th international workshop on cooperative and human aspects of software engineering (CHASE/14)*.

FIXME

Taschuk, M., and G. Wilson. 2017. Ten simple rules for making research software more robust. *PLoS Computational Biology*. 13. doi:10.1371/journal.pcbi.1005412.

A short guide to making research software usable by other people.

Udell, J. 2011. Seven ways to think like the web.

FIXME

Waldo, J. et al. 1994. A note on distributed computing. *IEEE Micro*.

FIXME

Whitehead, A.N. 1958. *An introduction to mathematics*. Oxford University Press.

FIXME

Wilson, G. et al. 2014. Best practices for scientific computing. *PLoS Biology*. 12. doi:10.1371/journal.pbio.1001745.

Outlines what a mature research software project should look like.

Wilson, G. et al. 2017. Good enough practices in scientific computing. *PLoS Computational Biology*. 13. doi:10.1371/journal.pcbi.1005510.

Outlines what a “good enough” research software project should look like.

Wing, J.M. 2006. Computational thinking. *Communications of the ACM*. 49. doi:10.1145/1118178.1118215.

First use of the term “computational thinking”.

Xu, T. et al. 2015. Hey, you have given me too many knobs!: Understanding and dealing with over-designed configuration in system software. *In Proc. 10th*

joint meeting on foundations of software engineering (FSE'2015). ACM Press. Examines the over-abundance of configuration options in software.

Yenni, G.M. et al. 2019. Developing a modern data workflow for regularly updated data. *PLOS Biology*. 17:e3000125. doi:10.1371/journal.pbio.3000125. FIXME

Yuan, D. et al. 2014. Simple testing can prevent most critical failures: An analysis of production failures in distributed data-intensive systems. *In* OSDI. Most failures in big data systems can be prevented by testing the error-handling code.