

# R for Excel Users

*Julie Lowndes & Allison Horst*

*2020-01-09*



# Contents

<b>1</b>	<b>Welcome</b>	<b>7</b>
1.1	Agenda . . . . .	7
1.2	Prerequisites . . . . .	7
<b>2</b>	<b>Overview</b>	<b>9</b>
2.1	Welcome! . . . . .	9
2.2	Why learn R if I know Excel? . . . . .	9
2.3	Guiding principles / recurring themes . . . . .	11
2.4	Resources . . . . .	12
<b>3</b>	<b>R &amp; RStudio, RMarkdown</b>	<b>13</b>
3.1	Summary . . . . .	13
3.2	RStudio Orientation . . . . .	13
3.3	Intro to RMarkdown . . . . .	15
3.4	Help pages . . . . .	20
3.5	Functions and Packages . . . . .	22
3.6	R Console . . . . .	24
3.7	Packages . . . . .	25
3.8	Clearing the environment . . . . .	27
3.9	Setup Git & GitHub . . . . .	28
3.10	Deep thoughts . . . . .	29
3.11	Efficiency Tips . . . . .	29
<b>4</b>	<b>GitHub</b>	<b>31</b>
4.1	Summary . . . . .	31
4.2	Objectives . . . . .	31
4.3	Resources . . . . .	32
4.4	Why should R users use Github? . . . . .	32
4.5	Create a repository on Github.com . . . . .	35
4.6	Create a gh-pages branch . . . . .	36
4.7	Clone your repository using RStudio . . . . .	37
4.8	Inspect your repository . . . . .	41
4.9	Add files to our local repo . . . . .	42
4.10	Sync from RStudio to GitHub . . . . .	44

4.11 Explore remote Github . . . . .	47
4.12 Create a new R Markdown file . . . . .	48
4.13 Explore your webpage . . . . .	49
4.14 Committing - how often? Tracking changes in your files . . . . .	50
4.15 Happy Git with R . . . . .	51
4.16 Efficiency Tips . . . . .	51
<b>5 Graphs with ggplot2</b> . . . . .	<b>53</b>
5.1 Summary . . . . .	53
5.2 Getting started - Create a new .Rmd, attach packages & get data . . . . .	55
5.3 Our first ggplot graph: Visitors to Channel Islands NP . . . . .	56
5.4 Intro to customizing ggplot graphs . . . . .	59
5.5 Mapping variables onto aesthetics . . . . .	61
5.6 ggplot2 complete themes . . . . .	63
5.7 Updating axis labels and titles . . . . .	64
5.8 Combining compatible geoms . . . . .	65
5.9 Multi-series ggplot graphs . . . . .	67
5.10 Faceting ggplot graphs . . . . .	69
5.11 Exporting a ggplot graph with <code>ggsave()</code> . . . . .	70
5.12 End ggplot session . . . . .	71
<b>6 dplyr and Pivot Tables</b> . . . . .	<b>73</b>
6.1 Summary . . . . .	73
6.2 Objectives . . . . .	73
6.3 Resources . . . . .	73
6.4 Pivot table overview . . . . .	74
6.5 RMarkdown setup . . . . .	74
6.6 Pivot table demo . . . . .	76
6.7 <code>group_by() %&gt;% summarize()</code> . . . . .	80
6.8 Oh no, our colleague sent the wrong data! . . . . .	85
6.9 <code>mutate()</code> . . . . .	89
6.10 <code>select()</code> . . . . .	90
6.11 Deep thoughts . . . . .	91
6.12 Efficiency Tips . . . . .	91
<b>7 Tidying</b> . . . . .	<b>93</b>
7.1 Summary . . . . .	93
7.2 Objectives . . . . .	93
7.3 Resources . . . . .	93
7.4 Set-up . . . . .	94
7.5 Wide-to-longer format with <code>tidyverse::pivot_longer()</code> . . . . .	94
7.6 Long-to-wider format with <code>tidyverse::pivot_wider()</code> . . . . .	96
7.7 Clean up column names with <code>janitor::clean_names()</code> . . . . .	97
7.8 Combine or separate columns with <code>tidyverse::unite()</code> and <code>tidyverse::separate()</code> . . . . .	98
7.9 Get just the numbers with <code>readr::parse_number()</code> . . . . .	102

7.10 Detecting a string pattern with <code>stringr::str_detect()</code>	103
7.11 Replacing a string with <code>stringr::str_replace()</code>	105
<b>8 Dplyr and vlookups</b>	<b>107</b>
8.1 Summary	107
8.2 Set-up: Create a new .Rmd, attach packages & get data	108
8.3 <code>dplyr::filter()</code> to conditionally subset by rows	109
8.4 Join data frames with <code>dplyr::*_join()</code>	113
8.5 A nice HTML table with <code>kable()</code> and <code>kableExtra</code>	118
8.6 Fun / kind of scary facts	118
8.7 Efficiency Tips	119
8.8 End <code>filter()</code> + <code>join()</code> section	119
<b>9 Collaborating</b>	<b>121</b>
9.1 Summary	121
9.2 Lesson	121
9.3 Troubleshooting: help you help yourself	121
9.4 Efficiency Tips	122
9.5 Additional thoughts	122
<b>10 Synthesis</b>	<b>123</b>
10.1 Summary	123
10.2 Objectives	123
10.3 Resources	124
10.4 Set-up	124
10.5 Attach packages, read in and explore the data	124
10.6 Some data cleaning to get salmon landings	125
10.7 Find grouped summary data	126
10.8 Make a graph of US commercial fisheries value by species over time with <code>ggplot2</code>	126
10.9 Export your salmon value graph with <code>ggsave</code>	127
10.10 2015 commercial fisheries value by state	127
10.11 Making a nice HTML table	128
10.12 Sync with GitHub remote	128



# Chapter 1

## Welcome

Hello! This is a workshop taught by Julie Stewart Lowndes and Allison Horst at the RStudio Conference: January 27-28 in San Francisco, California.

We are environmental scientists who use and teach R in our daily work. We both work at the University of California Santa Barbara: Julie is based at the National Center for Ecological Analysis and Synthesis as part of the Ocean Health Index team and leads Openscapes, and Allison is based at the Bren School of Environmental Science and Management as a lecturer of data science & statistics — and is also an Artist in Residence at RStudio.

### 1.1 Agenda

Time	Day 1	Day 2
9-10:30 break	Motivation, R & RStudio, RMarkdown (JL)	Tidying data (AH)
11-12:30 lunch	Intro to GitHub (JL)	dplyr & VLOOKUPs (AH)
13:30-15:00 break	ggplot2 & Charts (AH)	Collaborating in #rstudio (JL)
15:30-17:00	dplyr & Pivot Tables (JL)	Synthesis (AH)

### 1.2 Prerequisites

Before the training, please make sure you have done the following:

1. Download and install **up-to-date versions** of:

- R: <https://cloud.r-project.org>
  - RStudio: <http://www.rstudio.com/download>
2. Install the Tidyverse
  3. Create a an account:
    - <https://github.com>
1. Get comfortable: if you're not in a physical workshop, be set up with two screens if possible. You will be following along in RStudio on your own computer while also watching a virtual training or following this tutorial on your own.

# Chapter 2

## Overview

TODO: add Star Wars illustrations & Wickham R4DS illustration (as slides?)

### 2.1 Welcome!

In this workshop you will learn hands-on how to begin to interoperate between Excel and R. But this workshop is not only about learning R; we will learn R using additional software: RStudio and GitHub. These tools will help us develop good habits for working in a reproducible and collaborative way — critical attributes of the modern analyst.

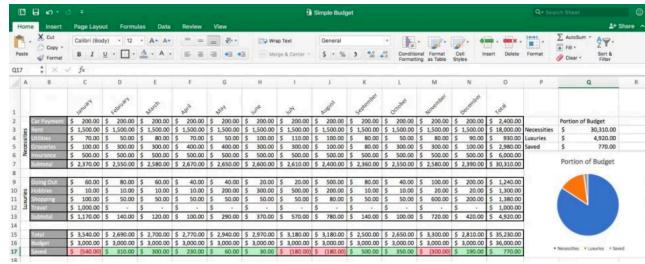
It's going to be fun and empowering!

### 2.2 Why learn R if I know Excel?

Excel is a widely used and powerful tool for working with data, and it is great for a lot of things. This is convenient and familiar; most of us have had their first experiences with data through Excel or other spreadsheet programs. As Jenny Bryan has said, “Excel is how we learn that we love data analysis”.

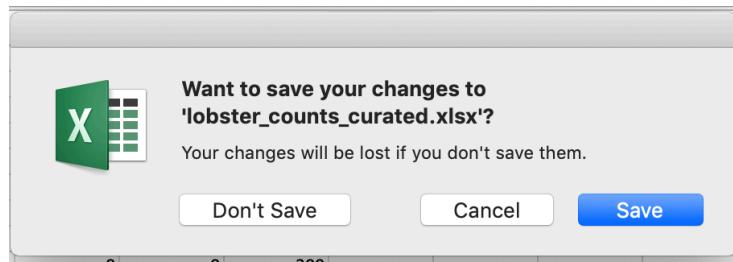
Excel is great for data entry. Can also be good for looking at data and feeling like you can touch it, and creating quick exploratory figures.

Excel can also become problematic with extending to analyses. This is because there aren't firm lines between what is data and what is analyses. For example, in this sheet:



This makes the analytical steps taken are not readily apparent, nor easy to reproduce. Have you ever done forensics on an Excel sheet, trying to understand what happened between columns or sheets? Maybe it was even your own Excel file from the (recent) past.

This also makes them pretty brittle/sensitive to minor changes. Has seeing this ever given you a feeling of horror:



So while it is great how easily you can update different fields and add analytical steps in an Excel sheet, it can also be a bit hard to handle, particularly as projects get more complicated.

So, as automation, reproducibility, collaboration, and frequent reporting become increasingly expected in data analysis, a good option for Excel users is to extend their workflows with R.

### 2.2.1 What to expect

This is going to be a fun workshop.

This workshop will give you hands-on experience and confidence with R, and how to interoperate between Excel and R — it is not about wholesale replacing everything you do in Excel into R. We will learn technical skills that you can incrementally incorporate into your existing workflows. But a big part of interfacing between Excel and R is not only skillsets, it is mindsets. It is the mindset about how we think about data. How we shape data and organize data and analyze data. And how what we do now can make our analytical life better in the future.

**A modern R user has a workflow framed around collaboration**, and uses an ecosystem of tools and practices. We will be learning three main things all at the same time:

1. coding with best practices (R/RStudio/tidyverse)
2. collaborative bookkeeping (Git/GitHub)
3. reporting and publishing (RMarkdown/GitHub)

**R users keep raw data separate from their analyses**, which means having data in one file and written computational commands saved as a separate file. We also embrace the concept of “**tidy data**”, where the data has a rectangular shape and each column is a variable and each row is an observation. Tidy data is a way of life.

region	state	code	park_name	type	visitors	year
PW	CA	CHIS	Channel Islands National Park	National Park	1200	1963
PW	CA	CHIS	Channel Islands National Park	National Park	1500	1964
PW	CA	CHIS	Channel Islands National Park	National Park	1600	1965
PW	CA	CHIS	Channel Islands National Park	National Park	300	1966
PW	CA	CHIS	Channel Islands National Park	National Park	15700	1967
PW	CA	CHIS	Channel Islands National Park	National Park	31000	1968
PW	CA	CHIS	Channel Islands National Park	National Park	33100	1969
PW	CA	CHIS	Channel Islands National Park	National Park	32000	1970

We are going to go through a lot in these two days and it's less important that you remember it all. More importantly, you'll have experience with it and confidence that you can do it. The main thing to take away is that there *are* good ways to work between R and Excel; we will teach you to expect that so you can find what you need and use it! A theme throughout is that tools exist and are being developed by real, and extraordinarily nice, people to meet you where you are and help you do what you need to do.

You are all welcome here, please be respectful of one another. Everyone in this workshop is coming from a different place with different experiences and expectations. But everyone will learn something new here, because there is so much innovation in the data science world. Instructors and helpers learn something new every time, from each other and from your questions. If you are already familiar with some of this material, focus on how we teach, and how you might teach it to others. Use these workshop materials not only as a reference in the future but also for talking points so you can communicate the importance of these tools to your communities. A big part of this training is not only for you to learn these skills, but for you to also teach others and increase the value and practice of open data science in science as a whole.

## 2.3 Guiding principles / recurring themes

**“Keep the raw data raw”** — A hard line separating raw data and analyses. In R, we have data in one file and written computational commands saved as a separate file.

**Scripted analyses** — We write analytical logic in code (rather than clicks) so that can be understood, rerun, and built upon.

**Learn from data that are not your own** — We aren't using your data in this workshop, but you will see similarities and patterns, and you'll see that these tools and practices apply to your work.

**Think ahead for Future You, Future Us.** Help make lives easier — first and foremost your own. Create breadcrumbs for yourselves and others: document and share your work.

## 2.4 Resources

R is not only a language, it is an active community of developers, users, and educators (often these traits are in each person). This workshop and book based on many excellent materials created by other members in the R community, who share their work freely to help others learn. Using community materials is how WE learned R, and each chapter of the book will have Resources listed for further reading into the topics we discuss. And, when there is no better way to explain something (ahem Jenny Bryan), we will quote or reference that work directly.

- What They Forgot to Teach You About R — Jenny Bryan & Jim Hester
- Stat545 — Jenny Bryan & Stat545 TAs
- Where do Things Live in R? REX Analytics
- 
- Spreadsheet Drama (Episode 9) — Not So Standard Deviations with Roger Peng & Hilary Parker
- more to come!

# Chapter 3

## R & RStudio, RMarkdown

### 3.1 Summary

We'll learn RMarkdown, which helps you tell a story with your data analysis because you can write text alongside the code. We are actually learning two languages at once: R and Markdown.

#### 3.1.1 Objectives

In this lesson we will:

- get oriented to the RStudio interface
- explore RMarkdown.
- discuss RMarkdown files vs Console (running vs knitting)
- learn a few base R functions (`c()`)
- error messages and help pages
- discuss and install packages (`here()`) usethis
- intro pipe operator (`%>%`)
- configure Git (to prepare for next session)

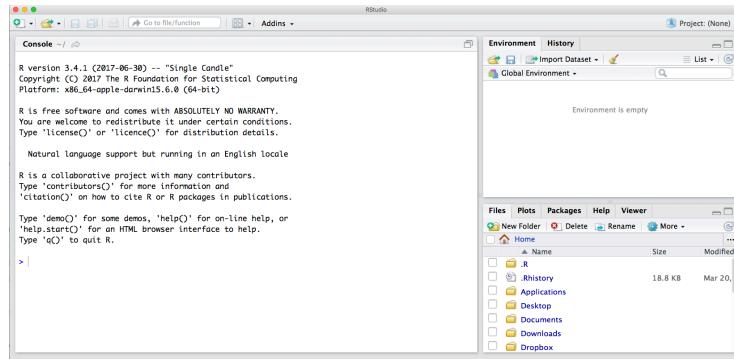
#### 3.1.2 Resources

- R for Excel Users by Gordon Shotwell (blog)

### 3.2 RStudio Orientation

Open RStudio for the first time.

Launch RStudio/R.



Notice the default panes:

- Console (entire left)
- Environment/History (tabbed in upper right)
- Files/Plots/Packages/Help (tabbed in lower right)

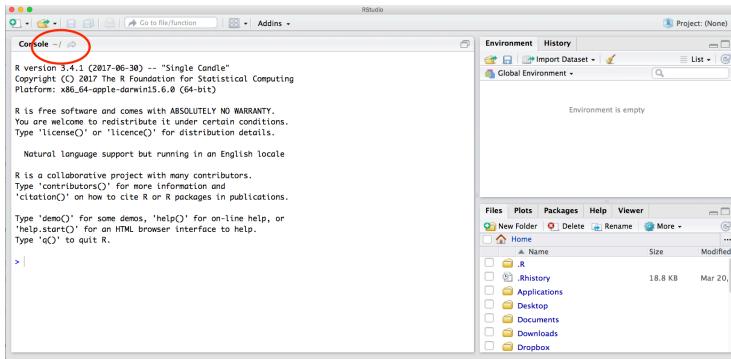
We won't click through this all immediately but we will become familiar with more of the options and capabilities throughout the next few days.

Something critical to know now is that you can make everything you see BIGGER by going to the navigation pane: View > Zoom In. Learn these keyboard shortcuts; being able to see what you're typing will help avoid typos & help us help you.

I think that **R is your airplane, and the RStudio IDE is your airport.** You are the pilot, and you use R to go places! With practice you'll gain skills and confidence; you can fly further distances and get through tricky situations. You will become an awesome pilot and can fly your plane anywhere. And the RStudio IDE provides support! Runways, communication, community, and other services that makes your life as a pilot much easier. It provides not only the infrastructure but a hub for the community that you can interact with.

An important first question: **where are we?**

If you've have opened RStudio for the first time, you'll be in your Home directory. This is noted by the `~` at the top of the console. You can see too that the Files pane in the lower right shows what is in the Home directory where you are. You can navigate around within that Files pane and explore, but note that you won't change where you are: even as you click through you'll still be Home: `~`.



We are going to have our first experience with R through RMarkdown, so let's do the following.

### 3.3 Intro to RMarkdown

An RMarkdown file is a plain text file that allow us to write code and text together, and when it is “knit”, the code will be evaluated and the text formatted so that it creates a reproducible report or document that is nice to read as a human.

This is really critical to reproducibility, and it also saves time. This document will recreate your figures for you in the same document where you are writing text. So no more doing analysis, saving a plot, pasting that plot into Word, redoing the analysis, re-saving, re-pasting, etc.

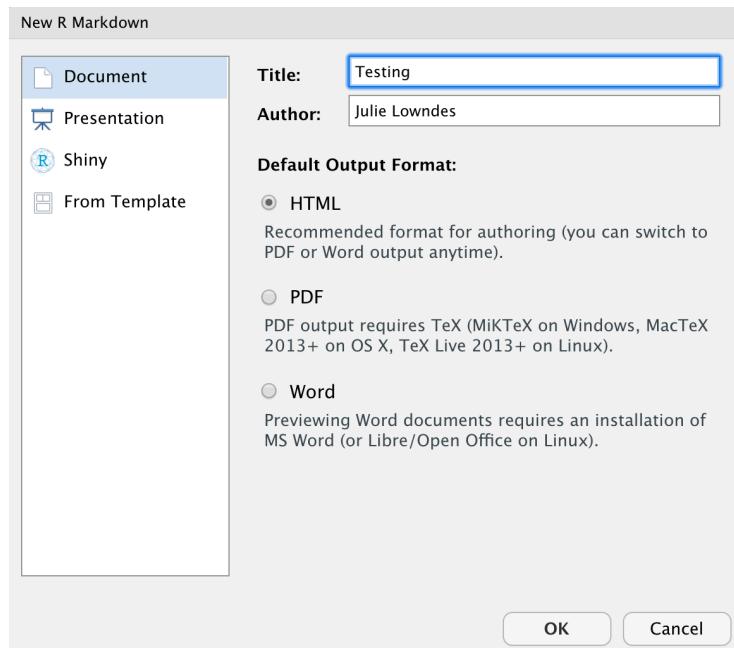
Let's experience this a bit ourselves and then we'll talk about it more.

#### 3.3.1 Create an RMarkdown file

Let's do this together:

File -> New File -> RMarkdown... (or alternatively you can click the green plus in the top left -> RMarkdown).

Let's title it “Testing” and write our name as author, then click OK with the recommended Default Output Format, which is HTML.



OK, first off: by opening a file, we are seeing the 4th pane of the RStudio console, which is essentially a text editor. This lets us organize our files within RStudio instead of having a bunch of different windows open.

Let's have a look at this file — it's not blank; there is some initial text is already provided for you. Let's have a high-level look through of it:

- The top part has the Title and Author we provided, as well as today's date and the output type as an HTML document like we selected above.
- There are white and grey sections. These are the 2 main languages that make up an RMarkdown file.
  - Grey sections are R code
  - White sections are Markdown text
- There is black and blue text.

```

1 ---  

2 title: "Testing"  

3 author: "Julie Lowndes"  

4 date: "11/14/2019"  

5 output: html_document  

6 ---  

7 |  

8 - ````{r setup, include=FALSE}  

9 knitr::opts_chunk$set(echo = TRUE)  

10 ````  

11  

12 - ## R Markdown  

13  

14 This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents.  

For more details on using R Markdown see <http://rmarkdown.rstudio.com>.  

15  

16 When you click the **Knit** button a document will be generated that includes both content as well as the output of any  

embedded R code chunks within the document. You can embed an R code chunk like this:  

17  

18 - ````{r cars}  

19 summary(cars)  

20 ````  

21  

22 - ## Including Plots  

23  

24 You can also embed plots, for example:  

25  

26 - ````{r pressure, echo=FALSE}  

27 plot(pressure)  

28 ````  

29  

30 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the  

plot.  

31

```

### 3.3.2 Knit your RMarkdown file

Let's go ahead and "Knit" by clicking the blue yarn at the top of the RMarkdown file. It's going to ask us to save first, I'll name mine "testing.Rmd". Note that this is by default going to save this file in your home directory `/~`. Since this is a testing document this is fine to save here; we will get more organized about where we save files very soon. Once you click Save, the knit process will be able to continue.

OK so how cool is this, we've just made an html file! This is a single webpage that we are viewing locally on our own computers. Knitting this RMarkdown document has rendered — we also say formatted — both the Markdown text (white) and the R code (grey), and the it also executed — we also say ran — the R code.

Let's have a look at them side-by-side:

The screenshot shows the RStudio interface with a knitted R Markdown document titled "Testing". The code editor on the left contains R code, including a summary of the "cars" dataset and a plot of "pressure" vs "temperature". The preview pane on the right displays the rendered Markdown, showing the table and the scatter plot.

```

1 #> 
2 #> title: "Testing"
3 #> author: "Julie Lowndes"
4 #> date: "11/14/2019"
5 #> output: html_document
6 #>
7 #>
8 #> # [r setup, include=FALSE]
9 knitr::opts_chunk$set(echo = TRUE)
10
11
12 #> # R Markdown
13
14 This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see http://rmarkdown.rstudio.com.
15
16 When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:
17
18<```{r cars}
19 summary(cars)
20```
21
22 #> # Including Plots
23
24 You can also embed plots, for example:
25
26<```{r pressure, echo=FALSE}
27 plot(pressure)
28```
29
30 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
279
280
281
282
283
284
285
286
287
288
289
289
290
291
292
293
294
295
296
297
298
299
299
300
301
302
303
304
305
306
307
308
309
309
310
311
312
313
314
315
316
317
318
319
319
320
321
322
323
324
325
326
327
328
329
329
330
331
332
333
334
335
336
337
338
339
339
340
341
342
343
344
345
346
347
348
349
349
350
351
352
353
354
355
356
357
358
359
359
360
361
362
363
364
365
366
367
368
369
369
370
371
372
373
374
375
376
377
378
379
379
380
381
382
383
384
385
386
387
388
389
389
390
391
392
393
394
395
396
397
398
399
399
400
401
402
403
404
405
406
407
408
409
409
410
411
412
413
414
415
416
417
418
419
419
420
421
422
423
424
425
426
427
428
429
429
430
431
432
433
434
435
436
437
438
439
439
440
441
442
443
444
445
446
447
448
449
449
450
451
452
453
454
455
456
457
458
459
459
460
461
462
463
464
465
466
467
468
469
469
470
471
472
473
474
475
476
477
478
479
479
480
481
482
483
484
485
486
487
488
489
489
490
491
492
493
494
495
496
497
498
499
499
500
501
502
503
504
505
506
507
508
509
509
510
511
512
513
514
515
516
517
518
519
519
520
521
522
523
524
525
526
527
528
529
529
530
531
532
533
534
535
536
537
538
539
539
540
541
542
543
544
545
546
547
548
549
549
550
551
552
553
554
555
556
557
558
559
559
560
561
562
563
564
565
566
567
568
569
569
570
571
572
573
574
575
576
577
578
579
579
580
581
582
583
584
585
586
587
588
589
589
590
591
592
593
594
595
596
597
598
599
599
600
601
602
603
604
605
606
607
608
609
609
610
611
612
613
614
615
616
617
618
618
619
620
621
622
623
624
625
626
627
628
629
629
630
631
632
633
634
635
636
637
638
639
639
640
641
642
643
644
645
646
647
648
649
649
650
651
652
653
654
655
656
657
658
659
659
660
661
662
663
664
665
666
667
668
669
669
670
671
672
673
674
675
676
677
678
679
679
680
681
682
683
684
685
686
687
688
689
689
690
691
692
693
694
695
696
697
697
698
699
699
700
701
702
703
704
705
706
707
708
709
709
710
711
712
713
714
715
716
717
717
718
719
719
720
721
722
723
724
725
726
727
728
729
729
730
731
732
733
734
735
736
737
738
739
739
740
741
742
743
744
745
746
747
748
749
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
769
770
771
772
773
774
775
776
777
778
778
779
779
780
781
782
783
784
785
786
787
787
788
789
789
790
791
792
793
794
795
796
796
797
798
798
799
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
848
849
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
877
878
878
879
879
880
881
882
883
884
885
886
887
887
888
889
889
890
891
892
893
894
895
895
896
896
897
897
898
898
899
899
900
901
902
903
904
905
905
906
906
907
907
908
908
909
909
910
910
911
911
912
912
913
913
914
914
915
915
916
916
917
917
918
918
919
919
920
920
921
921
922
922
923
923
924
924
925
925
926
926
927
927
928
928
929
929
930
930
931
931
932
932
933
933
934
934
935
935
936
936
937
937
938
938
939
939
940
940
941
941
942
942
943
943
944
944
945
945
946
946
947
947
948
948
949
949
950
950
951
951
952
952
953
953
954
954
955
955
956
956
957
957
958
958
959
959
960
960
961
961
962
962
963
963
964
964
965
965
966
966
967
967
968
968
969
969
970
970
971
971
972
972
973
973
974
974
975
975
976
976
977
977
978
978
979
979
980
980
981
981
982
982
983
983
984
984
985
985
986
986
987
987
988
988
989
989
990
990
991
991
992
992
993
993
994
994
995
995
996
996
997
997
998
998
999
999
1000
1000
1001
1001
1002
1002
1003
1003
1004
1004
1005
1005
1006
1006
1007
1007
1008
1008
1009
1009
1010
1010
1011
1011
1012
1012
1013
1013
1014
1014
1015
1015
1016
1016
1017
1017
1018
1018
1019
1019
1020
1020
1021
1021
1022
1022
1023
1023
1024
1024
1025
1025
1026
1026
1027
1027
1028
1028
1029
1029
1030
1030
1031
1031
1032
1032
1033
1033
1034
1034
1035
1035
1036
1036
1037
1037
1038
1038
1039
1039
1040
1040
1041
1041
1042
1042
1043
1043
1044
1044
1045
1045
1046
1046
1047
1047
1048
1048
1049
1049
1050
1050
1051
1051
1052
1052
1053
1053
1054
1054
1055
1055
1056
1056
1057
1057
1058
1058
1059
1059
1060
1060
1061
1061
1062
1062
1063
1063
1064
1064
1065
1065
1066
1066
1067
1067
1068
1068
1069
1069
1070
1070
1071
1071
1072
1072
1073
1073
1074
1074
1075
1075
1076
1076
1077
1077
1078
1078
1079
1079
1080
1080
1081
1081
1082
1082
1083
1083
1084
1084
1085
1085
1086
1086
1087
1087
1088
1088
1089
1089
1090
1090
1091
1091
1092
1092
1093
1093
1094
1094
1095
1095
1096
1096
1097
1097
1098
1098
1099
1099
1100
1100
1101
1101
1102
1102
1103
1103
1104
1104
1105
1105
1106
1106
1107
1107
1108
1108
1109
1109
1110
1110
1111
1111
1112
1112
1113
1113
1114
1114
1115
1115
1116
1116
1117
1117
1118
1118
1119
1119
1120
1120
1121
1121
1122
1122
1123
1123
1124
1124
1125
1125
1126
1126
1127
1127
1128
1128
1129
1129
1130
1130
1131
1131
1132
1132
1133
1133
1134
1134
1135
1135
1136
1136
1137
1137
1138
1138
1139
1139
1140
1140
1141
1141
1142
1142
1143
1143
1144
1144
1145
1145
1146
1146
1147
1147
1148
1148
1149
1149
1150
1150
1151
1151
1152
1152
1153
1153
1154
1154
1155
1155
1156
1156
1157
1157
1158
1158
1159
1159
1160
1160
1161
1161
1162
1162
1163
1163
1164
1164
1165
1165
1166
1166
1167
1167
1168
1168
1169
1169
1170
1170
1171
1171
1172
1172
1173
1173
1174
1174
1175
1175
1176
1176
1177
1177
1178
1178
1179
1179
1180
1180
1181
1181
1182
1182
1183
1183
1184
1184
1185
1185
1186
1186
1187
1187
1188
1188
1189
1189
1190
1190
1191
1191
1192
1192
1193
1193
1194
1194
1195
1195
1196
1196
1197
1197
1198
1198
1199
1199
1200
1200
1201
1201
1202
1202
1203
1203
1204
1204
1205
1205
1206
1206
1207
1207
1208
1208
1209
1209
1210
1210
1211
1211
1212
1212
1213
1213
1214
1214
1215
1215
1216
1216
1217
1217
1218
1218
1219
1219
1220
1220
1221
1221
1222
1222
1223
1223
1224
1224
1225
1225
1226
1226
1227
1227
1228
1228
1229
1229
1230
1230
1231
1231
1232
1232
1233
1233
1234
1234
1235
1235
1236
1236
1237
1237
1238
1238
1239
1239
1240
1240
1241
1241
1242
1242
1243
1243
1244
1244
1245
1245
1246
1246
1247
1247
1248
1248
1249
1249
1250
1250
1251
1251
1252
1252
1253
1253
1254
1254
1255
1255
1256
1256
1257
1257
1258
1258
1259
1259
1260
1260
1261
1261
1262
1262
1263
1263
1264
1264
1265
1265
1266
1266
1267
1267
1268
1268
1269
1269
1270
1270
1271
1271
1272
1272
1273
1273
1274
1274
1275
1275
1276
1276
1277
1277
1278
1278
1279
1279
1280
1280
1281
1281
1282
1282
1283
1283
1284
1284
1285
1285
1286
1286
1287
1287
1288
1288
1289
1289
1290
1290
1291
1291
1292
1292
1293
1293
1294
1294
1295
1295
1296
1296
1297
1297
1298
1298
1299
1299
1300
1300
1301
1301
1302
1302
1303
1303
1304
1304
1305
1305
1306
1306
1307
1307
1308
1308
1309
1309
1310
1310
1311
1311
1312
1312
1313
1313
1314
1314
1315
1315
1316
1316
1317
1317
1318
1318
1319
1319
1320
1320
1321
1321
1322
1322
1323
1323
1324
1324
1325
1325
1326
1326
1327
1327
1328
1328
1329
1329
1330
1330
1331
1331
1332
1332
1333
1333
1334
1334
1335
1335
1336
1336
1337
1337
1338
1338
1339
1339
1340
1340
1341
1341
1342
1342
1343
1343
1344
1344
1345
1345
1346
1346
1347
1347
1348
1348
1349
1349
1350
1350
1351
1351
1352
1352
1353
1353
1354
1354
1355
1355
1356
1356
1357
1357
1358
1358
1359
1359
1360
1360
1361
1361
1362
1362
1363
1363
1364
1364
1365
1365
1366
1366
1367
1367
1368
1368
1369
1369
1370
1370
1371
1371
1372
1372
1373
1373
1374
1374
1375
1375
1376
1376
1377
1377
1378
1378
1379
1379
1380
1380
1381
1381
1382
1382
1383
1383
1384
1384
1385
1385
1386
1386
1387
1387
1388
1388
1389
1389
1390
1390
1391
1391
1392
1392
1393
1393
1394
1394
1395
1395
1396
1396
1397
1397
1398
1398
1399
1399
1400
1400
1401
1401
1402
1402
1403
1403
1404
1404
1405
1405
1406
1406
1407
1407
1408
1408
1409
1409
1410
1410
1411
1411
1412
1412
1413
1413
1414
1414
1415
1415
1416
1416
1417
1417
1418
1418
1419
1419
1420
1420
1421
1421
1422
1422
1423
1423
1424
1424
1425
1425
1426
1426
1427
1427
1428
1428
1429
1429
1430
1430
1431
1431
1432
1432
1433
1433
1434
1434
1435
1435
1436
1436
1437
1437
1438
1438
1439
1439
1440
1440
1441
1441
1442
1442
1443
1443
1444
1444
1445
1445
1446
1446
1447
1447
1448
1448
1449
1449
1450
1450
1451
1451
1452
1452
1453
1453
1454
1454
1455
1455
1456
1456
1457
1457
1458
1458
1459
1459
1460
1460
1461
1461
1462
1462
1463
1463
1464
1464
1465
1465
1466
1466
1467
1467
1468
1468
1469
1469
1470
1470
1471
1471
1472
1472
1473
1473
1474
1474
1475
1475
1476
1476
1477
1477
1478
1478
1479
1479
1480
1480
1481
1481
1482
1482
1483
1483
1484
1484
1485
1485
1486
1486
1487
1487
1488
1488
1489
1489
1490
1490
1491
1491
1492
1492
1493
1493
1494
1494
1495
1495
1496
1496
1497
1497
1498
1498
1499
1499
1500
1500
1501
1501
1502
1502
1503
1503
1504
1504
1505
1505
1506
1506
1507
1507
1508
1508
1509
1509
1510
1510
1511
1511
1512
1512
1513
1513
1514
1514
1515
1515
1516
1516
1517
1517
1518
1518
1519
1519
1520
1520
1521
1521
1522
1522
1523
1523
1524
1524
1525
1525
1526
1526
1527
1527
1528
1528
1529
1529
1530
1530
1531
1531
1532
1532
1533
1533
1534
1534
1535
1535
1536
1536
1537
1537
1538
1538
1539
1539
1540
1540
1541
1541
1542
1542
1543
1543
1544
1544
1545
1545
1546
1546
1547
1547
1548
1548
1549
1549
1550
1550
1551
1551
1552
1552
1553
1553
1554
1554
1555
1555
1556
1556
1557
1557
1558
1558
1559
1559
1560
1560
1561
1561
1562
1562
1563
1563
1564
1564
1565
1565
1566
1566
1567
1567
1568
1568
1569
1569
1570
1570
1571
1571
1572
1572
1573
1573
1574
1574
1575
1575
1576
1576
1577
1577
1578
1578
1579
1579
1580
1580
1581
1581
1582
1582
1583
1583
1584
1584
1585
1585
1586
1586
1587
1587
1588
1588
1589
1589
1590
1590
1591
1591
1592
1592
1593
1593
1594
1594
1595
1595
1596
1596
1597
1597
1598
1598
1599
1599
1600
1600
1601
1601
1602
1602
1603
1603
1604
1604
1605
1605
1606
1606
1607
1607
1608
1608
1609
1609
1610
1610
1611
1611
1612
1612
1613
1613
1614
1614
1615
1615
1616
1616
1617
1617
1618
1618
1619
1619
1620
1620
1621
1621
1622
1622
1623
1623
1624
1624
1625
1625
1626
1626
1627
1627
1628
1628
1629
1629
1630
1630
1631
1631
1632
1632
1633
1633
1634
1634
1635
1635
1636

```

Each of them are start by 3 backticks and `{r label}`, which signify that there will be R code inside, and they are each given a unique name. Anything inside the brackets (`{ }`) is instructions for RMarkdown about that code to run. For example:

- the first chunk says `include=FALSE`, and we don't see it included in the HTML document.
- the second chunk has no additional instructions, and in the HTML document we see the code and the evaluation of that code (a summary table)
- the third chunk says `echo=FALSE`, and in the HTML document we do not see the code echoed, we only see the plot when the code is executed.

There are many more options available and we will explore more as we go.

### 3.3.5.1 Naming code chunks (Deep thought?)

All three chunks say `r` as the language, and have a label (`setup`, `cars`, `pressure`). This is to help us navigate between them and keep them organized. They are optional, but will become powerful as you become a powerful R user. But if you label your code chunks, you must have unique labels. Otherwise you will see an error when you try to knit:

```
processing file: Untitled.Rmd
Error in parse_block(g[-1], g[1], params.src) : duplicate label 'cars'
Calls: <Anonymous> ... process_file -> split_file -> lapply -> FUN -> parse_block
Execution halted
```

In this case, you read the error message. Not everything immediately looks like something I would know anything about, but pressing on still allows me to identify the problem: “duplicate label ‘cars’”.

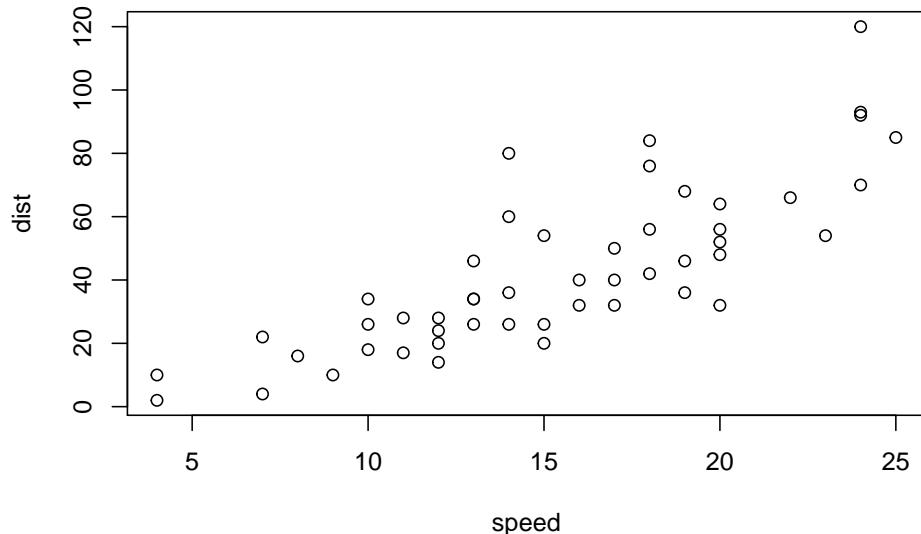
### 3.3.5.2 New code chunks

We can create a new chunk in your RMarkdown first in one of these ways:

- click “Insert > R” at the top of the editor pane
- type by hand “`{r}`”
- if you haven’t deleted a chunk that came with the new file, edit that one

Now, let’s write some code in R. Let’s say we want to plot the cars data. I’m going to press enter to to add some extra carriage returns because sometimes I find it easier to look at my code when there is a bit more space, and R lets you use as much whitespace as you would like.

```
plot(cars)
```



We can knit this and see the plot of cars. This is the same data that we made the summary table with.

Troubleshooting: If you could not successfully knit your document without error, Look at your code again. Do you have both open ( and close ) parentheses? Are your code chunk fences correct?

What is `plot` anyways? It is a function. Let's talk about that next.

### 3.4 Help pages

*TODO - complete*

Like Excel, some of the biggest power in R is that there are built-in functions that you can use in your analyses (and, as we'll see, R users can easily create and share functions, and it is this open source developer and contributor community that makes R so awesome).

R has a mind-blowing collection of built-in functions that are used with the same syntax: function name with parentheses around what the function needs to do what it is supposed to do.

TODO: Add Horst illustration

`function_name(argument1 = value1, argument2 = value2, ...)`. When you see this syntax, we say we are “calling the function”.

A good way to learn about functions is to look at the help pages. Let's navigate over to the Help tab and type `plot`. All help pages will have the same format, here is how I look at it:

The help page tells the name of the package in the top left, and broken down into sections:

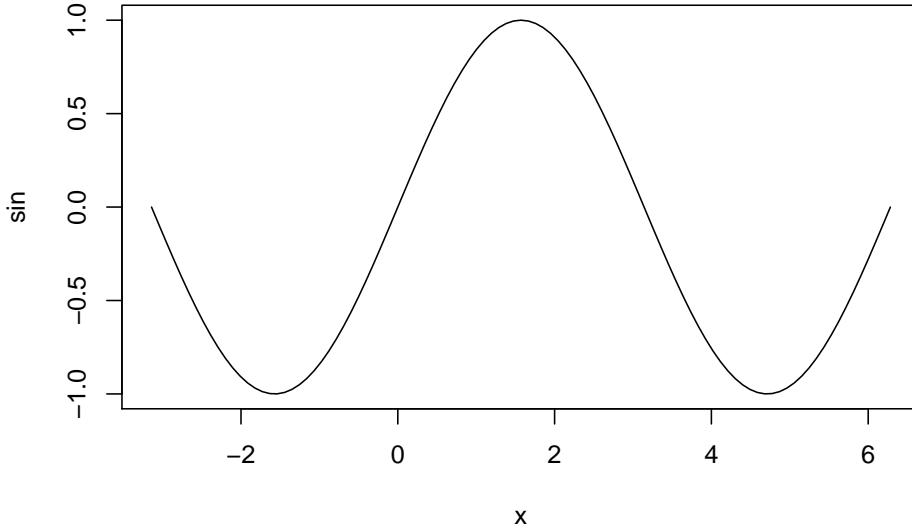
- Description: An extended description of what the function does.
- Usage: The arguments of the function and their default values.
- Arguments: An explanation of the data each argument is expecting.
- Details: Any important details to be aware of.
- Value: The data the function returns.
- See Also: Any related functions you might find useful.
- Examples: Some examples for how to use the function.

When I look at a help page, I start with the description, which might be too in-the-weeds for the level of understanding I need at the offset. For the `sum` page, it is pretty straight-forward and lets me know that yup, this is the function I want.

I next look at the usage and arguments, which give me a more concrete view into what the function does. This syntax looks a bit cryptic but what it means is that you use it by writing `sum`, and then passing whatever you want to it in terms of data: that is what the “`...`” means. And the “`na.rm=FALSE`” means that by default, it will not remove NAs (I read this as: “remove NAs? FALSE!”)

Then, I usually scroll down to the bottom to the examples. This is where I can actually see how the function is used, and I can also paste those examples into the Console to see their output. Let’s try it:

```
plot(sin, -pi, 2*pi)
```



#TODO: transition

Let’s try using `c()` which combines values together and, while we’re at it, demo more helpful features of RStudio.

So let's create a new R code chunk. And we'll write:

```
c(1, 5:9)
```

```
## [1] 1 5 6 7 8 9
```

So you can see that this combines these values all into the same place, which is called a vector here. But let's learn more about what this does.

### 3.4.1 Help pages NOPE

When I press enter/return, it will open up a help page in the bottom right pane. Help pages vary in detail I find some easier to digest than others. But they all have the same structure, which is helpful to know.

```
seq(from = 1, to = 10) # same as seq(1, 10); R assumes by position
```

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

```
seq(from = 1, to = 10, by = 2)
```

```
## [1] 1 3 5 7 9
```

The above also demonstrates something about how R resolves function arguments. You can always specify in `name = value` form. But if you do not, R attempts to resolve by position. So above, it is assumed that we want a sequence `from = 1` that goes `to = 10`. Since we didn't specify step size, the default value of `by` in the function definition is used, which ends up being 1 in this case. For functions I call often, I might use this resolve by position for the first argument or maybe the first two. After that, I always use `name = value`.

The examples from the help pages can be copy-pasted into the console for you to understand what's going on.

Not all functions have (or require) arguments:

```
date()
```

```
## [1] "Thu Jan 9 06:00:35 2020"
```

## 3.5 Functions and Packages

### 3.5.1 R code in the Console – where does this go

Now, hitting return does not execute this command; remember, it's a text file in the text editor, it's not associated with the R engine. To execute it, we need to get what we typed in the the R chunk (the grey R code) down into the console. How do we do it? There are several ways (let's do each of them):

1. copy-paste this line into the console.
2. select the line (or simply put the cursor there), and click ‘Run’. This is available from
  - a. the bar above the file (green arrow)
  - b. the menu bar: Code > Run Selected Line(s)
  - c. keyboard shortcut: command-return
3. click the green arrow at the right of the code chunk

And when we do this, we see the `mean_dist` object appear in the Environment pane.

Cool tip: doesn’t have to be only R, other languages supported.

Learn more: <http://rmarkdown.rstudio.com/>

### 3.5.1.1 Activity

1. In Markdown write some italic text, make a numbered list, and add a few subheaders. Use the Markdown Quick Reference (in the menu bar: Help > Markdown Quick Reference).
2. Reknit your html file.

### 3.5.1.2 Restart R

To end our work from this session, save, knit, and then restart R (Go to the top menus: Session > Restart R.)

Notice that now with a clean workspace, if I knit my document instead of sending code to the Console, my objects (like `mean_dist`) don’t show up in my Environment. This is because R isn’t actually running this in this R session, it is actually spinning up a clean session to knit my document. This is important for reproducible analyses because I don’t want the success of this analysis to be dependent on some weird setting I have on my computer that will make Future Me or Future Us not able to run or understand these important analyses. Having RMarkdown be self-contained in this way helps you develop good habits for reproducibility.

## 3.5.2 What is RMarkdown? (1-minute video)

Let’s watch this to demonstrate all the amazing things you can now do:

What is RMarkdown?

## 3.6 R Console

OK let's go into the Console, where we interact with the live R process.

We can do math:

```
52*40
365/12
```

But like Excel, the power comes not from doing small operations by hand (like  $8*22.3$ ), it's by being able to operate on whole suites of numbers and datasets. In Excel, data are stored in the spreadsheet. In R, they are stored in objects, which are often vectors or dataframes. They are rectangular.

### 3.6.1 Viewing data in R

Let's have a look at some data in R. Unlike Excel, R comes out-of-the-box with several built-in data sets that we can look at and work with.

One of these datasets is called `cars`. If I write this in the Console, it will print the data in the console.

```
cars
```

This returns data. To me this is not super interesting data, but I can appreciate that there are different variables listed as column headers and then numeric values for each type of row observation. (Unfortunately I don't know if these are different cars or trials or conditions but we won't focus on that for now).

I can also use RStudio's Viewer to see this in a more familiar-looking format. Let's type this — and make sure it's a capital V and open and closed parentheses:

```
View(cars)
```

This opens the fourth pane of the RStudio IDE; when you work in R you will have all four panes open so this will become a very comforting setup for you.

**Aside** The basic R data structure is a vector. You can think of a vector like a column in an Excel spreadsheet with the limitation that all the data in that vector must be of the same type. If it is a character vector, every element must be a character; if it is a logical vector, every element must be TRUE or FALSE; if it's numeric you can trust that every element is a number. There's no such constraint in Excel: you might have a column which has a bunch of numbers, but then some explanatory text intermingled with the numbers. This isn't allowed in R. - Gordon Shotwell

In the Viewer I can do things like filter or sort. This does not do anything to the actual data, it just changes how you are viewing the data. So even as I explore it, I am not editing or manipulating the data.

## 3.7 Packages

So far we've been using a couple functions from base R, such as `sum()` and `date()`. But, one of the amazing things about R is that a vast user community is always creating new functions and packages that expand R's capabilities. In R, the fundamental unit of shareable code is the package. A package bundles together code, data, documentation, and tests, and is easy to share with others. They increase the power of R by improving existing base R functionalities, or by adding new ones.

The traditional place to download packages is from CRAN, the Comprehensive R Archive Network, which is where you downloaded R. You can also install packages from GitHub, which we'll do tomorrow.

You don't need to go to CRAN's website to install packages, this can be accomplished within R using the command `install.packages("package-name-in-quotes")`. Let's install a small, fun package `praise`. You need to use quotes around the package name.

Do this in the Console instead of in your RMarkdown file because we don't want this to load every time:

```
install.packages("praise")
```

Now we've installed the package, but we need to tell R that we are going to use the functions within the `praise` package. We do this by using the function `library()`.

In my mind, this is analogous to needing to wire your house for electricity: this is something you do once; this is `install.packages`. But then you need to turn on the lights each time you need them (R Session).

```
library(praise)
```

Now that we've loaded the `praise` package, we can use the single function in the package, `praise()`, which returns a randomized praise to make you feel better.

```
praise()
```

```
## [1] "You are bee's knees!"
```

### 3.7.1 Assigning objects with <-

This might be a really important value that we want to have on hand. We can save this as its own object.

This is a big difference with Excel, where you usually identify data by its location on the grid, like `$A1:D$20`. (You can do this with Excel by naming ranges of

cells, but many people don't do this.)

Data can be a variety of formats, like numeric and text.

We do this by writing the name along with the assignment operator <-

```
sum_dist <- sum(cars$dist)
```

And we can execute it. In my head I hear "sum\_dist gets 2149".

Object names can be whatever you want, although they cannot start with a digit and cannot contain certain other characters such as a comma or a space. Different folks have different conventions; you will be wise to adopt a convention for demarcating words in names.

```
# i_use_snake_case
# other.people.use.periods
# evenOthersUseCamelCase
```

Notice that as I start typing `sum_dist` in the Console, there will be options to auto-fill. This is RStudio helping you out, which is great because we all are prone to typos. I actually have to ignore the help to try to force a typo:

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

OK this is an error, but I didn't break R — error messages are your friends.

The first thing to do with an error message is read it. Yes it's in angry red text and it's unexpected — but most error messages are doing their best to help you solve the problem. And you'll get more familiar with they way they tell you. By saying "object 'sumdist' not found" alerts me immediately to the fact that this thing I think exists R doesn't think exists — so maybe it's a typo or not loaded?

### 3.7.2 Error messages are your friends

As Jenny Bryan says:

Implicit contract with the computer / scripting language: Computer will do tedious computation for you. In return, you will be completely precise in your instructions. Typos matter. Case matters. Pay attention to how you type.

Remember that this is a language, not unsimilar to English! There are times you aren't understood — it's going to happen. There are different ways this can happen. Sometimes you'll get an error. This is like someone saying 'What?' or 'Pardon'? Error messages can also be more useful, like when they say 'I didn't understand what you said, I was expecting you to say blah'. That is

a great type of error message. Error messages are your friend. Google them (copy-and-paste!) to figure out what they mean.

And also know that there are errors that can creep in more subtly, when you are giving information that is understood, but not in the way you meant. Like if I am telling a story about suspenders that my British friend hears but silently interprets in a very different way (true story). This can leave me thinking I've gotten something across that the listener (or R) might silently interpreted very differently. And as I continue telling my story you get more and more confused... Clear communication is critical when you code: write clean, well documented code and check your work as you go to minimize these circumstances!

**Shortcuts** You will make lots of assignments and the operator `<-` is a pain to type. Don't be lazy and use `=`, although it would work, because it will just sow confusion later. Instead, utilize **RStudio's keyboard shortcut: Alt + - (the minus sign)**. Notice that RStudio automagically surrounds `<-` with spaces, which demonstrates a useful code formatting practice. Code is miserable to read on a good day. Give your eyes a break and use spaces. RStudio offers many handy keyboard shortcuts. Also, Alt+Shift+K brings up a keyboard shortcut reference card.

My most common shortcuts include command-Z (undo), and combinations of arrow keys in combination with shift/option/command (moving quickly up, down, sideways, with or without highlighting).

## 3.8 Clearing the environment

Now look at the objects in your environment (workspace) – in the upper right pane. The workspace is where user-defined objects accumulate.

For reproducibility, it is critical that you delete your objects and restart your R session frequently. You don't want your whole analysis to only work in whatever way you've been working right now — you need it to work next week, after you upgrade your operating system, etc. Restarting your R session will help you identify and account for anything you need for your analysis.

We will keep coming back to this theme but let's restart our R session together: Go to the top menus: Session > Restart R.

Don't save the workspace!

So this is great, but if we were going to do any kind of real analysis, we need to be able to write it in a document rather than this command line prompt in the Console. Let's do something much more interesting and really start feeling its power.

## 3.9 Setup Git & GitHub

Before we break, we are going to set up Git and GitHub for our next lesson (and if you have trouble please stay a few extra minutes).

We're going to switch gears from R for a moment and set up Git and GitHub, which we will be using along with R and RStudio for the rest of the workshop. This set up is a one-time thing! You will only have to do this once per computer. We'll walk through this together.

1. We will use the `usethis` package to configure `git` with global commands, which means it will apply ‘globally’ to all files on your computer, rather than to a specific folder.

```
install.packages("usethis")
library(usethis)

use_git_config(user.name = "jules32", user.email = "jules32@example.org")
```

*BACKUP PLAN* If `usethis` fails, the following is the classic approach to configuring `git`. Open the Git Bash program (Windows) or the Terminal (Mac) and type the following:

```
# display your version of git
git --version

# replace USER with your Github user account
git config --global user.name USER

# replace NAME@EMAIL.EDU with the email you used to register with Github
git config --global user.email NAME@EMAIL.EDU

# list your config to confirm user.* variables set
git config --list
```

Not only have you just set up git as a one-time-only thing, you have just used the command line. We don't have time to learn much of the command line today, but you just successfully used it following explicit instructions, which is huge! There are great resources for learning the command line, check out this tutorial from SWC at UCSB.

### 3.9.1 Troubleshooting

If you have problems setting up git, please see the Troubleshooting section in Jenny Bryan's amazing HappyGitWithR.

### 3.9.1.1 New(ish) Error on a Mac

We've also seen the following errors from RStudio:

```
error key does not contain a section --global terminal
```

and

```
fatal: not in a git directory
```

To solve this, go to the Terminal and type: `which git`

Look at the filepath that is returned. Does it say anything to do with Apple?

-> If yes, then the Git you downloaded isn't installed, please redownload if necessary, and follow instructions to install.

-> If no, (in the example image, the filepath does not say anything with Apple) then proceed below:

In RStudio, navigate to: Tools > Global Options > Git/SVN.

Does the “**Git executable**” filepath match what the url in Terminal says?

If not, click the browse button and navigate there.

*Note:* on my laptop, even though I navigated to /usr/local/bin/git, it then automatically redirect because /usr/local/bin/git was an alias on my computer. That is fine. Click OK.

Quit RStudio.

Then relaunch RStudio.

Try syncing or cloning, and if that works and then you don't need to worry about typing into the Terminal, you're all done!

## 3.10 Deep thoughts

Comments! Organization (spacing, subsections, vertical structure, indentation, etc.)! Well-named variables! Also, well-named operations so analyses (`select(data, columnname)`) instead of `data[1:6,5]` and excel equivalent. (Ex with strings) Not so brittle/sensitive to minor changes.

## 3.11 Efficiency Tips

—>



# Chapter 4

# GitHub

TODO: no github folder, new screenshots, earlier emphasis on syncing steps

## 4.1 Summary

We will learn about version control using git and GitHub, and we will interface with this through RStudio. Why use version control? To save time when working with your most important collaborator: you.

## 4.2 Objectives

Today, we'll interface with GitHub from our local computers using RStudio. There are many other ways to interact with GitHub, including GitHub's Desktop App or the command line (here is Jenny Bryan's list of git clients), but today we are going to work from RStudio. You have the largest suite of options if you interface through the command line, but the most common things you'll do can be done through one of these other applications (i.e. RStudio and the GitHub Desktop App).

Here's what we'll do after we set up git on your computers:

1. create a repository on Github.com
2. clone locally using RStudio
3. learn the RStudio-GitHub workflow by syncing to Github.com: pull, stage, commit, push
4. explore github.com: files, commit history, file history
5. practice the RStudio-GitHub workflow by editing and adding files
6. practice R Markdown

git will track and version your files, GitHub stores this online and enables you to collaborate with others (and yourself). Although git and GitHub are two different things, distinct from each other, we can think of them as a bundle since we will always use them together. It also helped me to think of GitHub like Dropbox: you make folders that are ‘tracked’ and can be synced to the cloud. GitHub does this too, but you have to be more deliberate about when syncs are made. This is because GitHub saves these as different versions, with information about who contributed when, line-by-line. This makes collaboration easier, and it allows you to roll-back to different versions or contribute to others’ work.

### 4.3 Resources

These materials borrow from:

- Jenny Bryan’s lectures from STAT545 at UBC: The Shell
- Jenny Bryan’s Happy git with R tutorial
- Melanie Frazier’s GitHub Quickstart, GitHub Lesson at University of Queensland
- Ben Best’s Software Carpentry at UCSB

Today, we’ll only introduce the features and terminology that new R users need to learn to begin managing their projects.

### 4.4 Why should R users use Github?

1. Ends (or, nearly ends) the horror of keeping track of versions. Basically, we

<input type="checkbox"/> Name	Date modified	Type
Rscript_4_21_2016.R	5/1/2016 3:03 PM	R File
Rscript_4_22_2016a.R	5/1/2016 3:03 PM	R File
Rscript_4_22_2016b.R	5/1/2016 3:03 PM	R File
Rscript_4_24_2016.R	5/1/2016 3:03 PM	R File
Rscript_final.R	5/1/2016 3:03 PM	R File
Rscript_final_final.R	5/1/2016 3:03 PM	R File
Rscript_really_final.R	5/1/2016 3:03 PM	R File
Rscript_really_really_final.R	5/1/2016 3:03 PM	R File

get away from this:

When you open your repository, you only see the most recent version.

But, it's easy to compare versions, and you can easily revert to previous versions.

2. Improves collaborative efforts. Different researchers can work on the same files at the same time!
3. It is easy to share and distribute files through the Github website.
4. Your files are available anywhere, you just need internet connection!

#### 4.4.1 What are Git and Github?

- **Git** is a version control system that lets you track changes to files over time. These files can be any kind of file (eg .doc, .pdf, .xls), but free text differences are most easily visible (eg txt, csv, md).
- **Github** is a website for storing your git versioned files remotely. It has many nice features to be able to visualize differences between images, rendering & diffing map data files, render text data files, and track changes in text.

Github was developed for social coding (i.e., sort of like an open source Wikipedia for programmers). Consequently, much of the functionality and terminology of Github (e.g., branches and pull requests) isn't necessary for a new R user getting started.

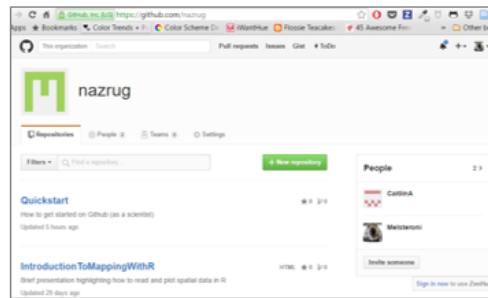
These concepts are more important for coders who want the entire coding community (and not just people working on the same project) to be able to suggest changes to their code. This isn't how most new R users will use Github.

To get the full functionality of Github, you will eventually want to learn other concepts. But, this can wait.

#### 4.4.2 Some Github terminology

- **User:** A Github account for you (e.g., jules32).
- **Organization:** The Github account for one or more users (e.g., datacar-pentry).
- **Repository:** A folder within the organization that includes files dedicated to a project.
- **Local Github:** Copies of Github files located on your computer.
- **Remote Github:** Github files located on the <https://github.com> website.
- **Clone:** Process of making a local copy of a remote Github repository. This only needs to be done once (unless you mess up your local copy).
- **Pull:** Copy changes on the remote Github repository to your local Github repository. This is useful if multiple people are making changes to a repository.
- **Push:** Save local changes to remote Github

## REMOTE (aka Github website)



**Clone** (i.e., copy)  
repository to your  
computer (a one  
time event)

Pull remote  
changes

Push local  
changes



**LOCAL**  
(aka your computer)

## 4.5 Create a repository on Github.com

First, go to your account on github.com and click “New repository”.

The screenshot shows a GitHub organization named "nazrug". The main navigation bar includes "Repositories", "People", "Teams", and "Settings". A large orange arrow points to the "+ New repository" button. Below the navigation, there are two repository cards: "IntroductionToMappingWithR" and "ManagingData". To the right, a sidebar lists "People" with profiles for "CaitlinA" and "Meisteroni", and a "Invite someone" button.

Choose a name. Call it whatever you want (the shorter the better), or follow me for convenience. I will call mine **r-workshop**.

Also, add a description, make it public, create a README file, and create your

The screenshot shows the "Create a new repository" form. A series of blue arrows with numbered steps guide the user through the process:

1. A simple name for the repository
2. A description
3. Public or private (private available with paid or education)
4. Select this
5. I ignore these
6. Click here!

The form fields include:

- Owner:** nazrug / Quickstart
- Description (optional):** How to get started on Github (as a scientist)
- Visibility:** Public (selected) vs Private
- Initialize this repository with a README:** Checked
- Add .gitignore:** None
- Add a license:** None
- Create repository:** Green button at the bottom

repo!

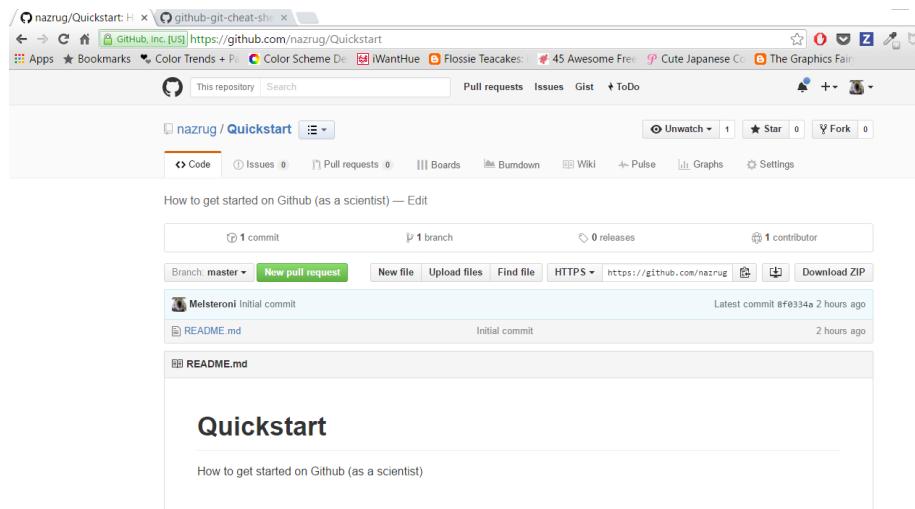
The *Add .gitignore* option adds a document where you can identify files or file-

types you want Github to ignore. These files will stay in on the local Github folder (the one on your computer), but will not be uploaded onto the web version of Github.

The *Add a license* option adds a license that describes how other people can use your Github files (e.g., open source, but no one can profit from them, etc.). We won't worry about this today.

Check out our new repository!

Notice how the README.md file we created is automatically displayed at the bottom. The .md means that it is Markdown (remember how .Rmd was RMarkdown?) so the formatting we learned in the last lesson apply here.



## 4.6 Create a gh-pages branch

We aren't going to talk about branches very much, but they are a powerful feature of git/GitHub. I think of it as creating a copy of your work that becomes a parallel universe that you can modify safely because it's not affecting your original work. And then you can choose to merge the universes back together if and when you want. By default, when you create a new repo you begin with one branch, and it is named **master**. When you create new branches, you can name them whatever you want. However, if you name one **gh-pages** (all lowercase, with a - and no spaces), this will let you create a website. And that's our plan. So, let's do this to create a **gh-pages** branch:

On the homepage for your repo on GitHub.com, click the button that says "Branch:master". Here, you can switch to another branch (right now there aren't any others besides **master**), or create one by typing a new name.

Let's type `gh-pages`.

Let's also change `gh-pages` to the default branch and delete the master branch: this will be a one-time-only thing that we do here:

First click to control branches:

And then click to change the default branch to `gh-pages`. I like to then delete the `master` branch when it has the little red trash can next to it. It will make you confirm that you really want to delete it, which I do!

**From here, you will work locally (on your computer).**

## 4.7 Clone your repository using RStudio

We'll start off by cloning to our local computer using RStudio. We are going to be cloning a copy of our Remote repository on Github.com to our local computers. Unlike downloading, cloning keeps all the version control and user information bundled with the files.

**Step 0:** Create your `github` folder

This is really important! We need to be organized and deliberate about where we want to keep all of our GitHub repositories (since this is the first of many in your career).

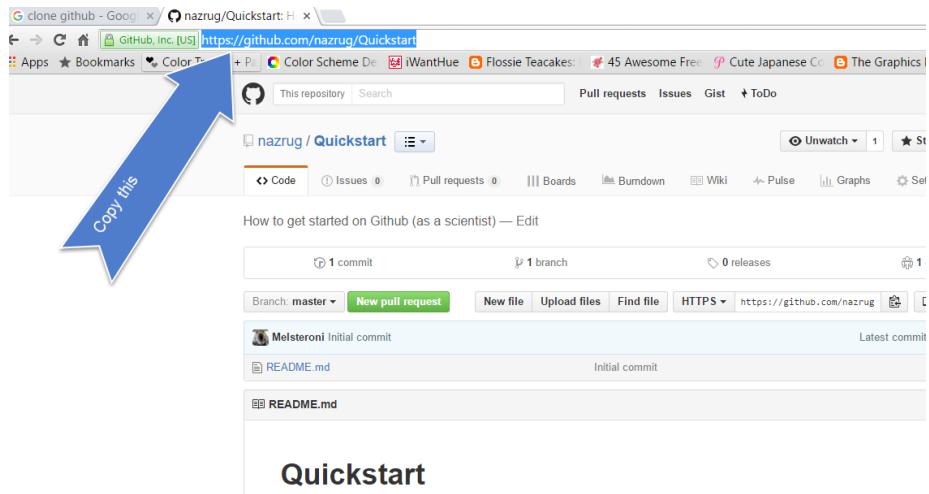
Let's all make a folder called `github` (all lowercase!) in our home directories. So it will look like this:

- Windows: `Users\[User]\Documents\github\`
- Mac: `Users/[User]/github/`

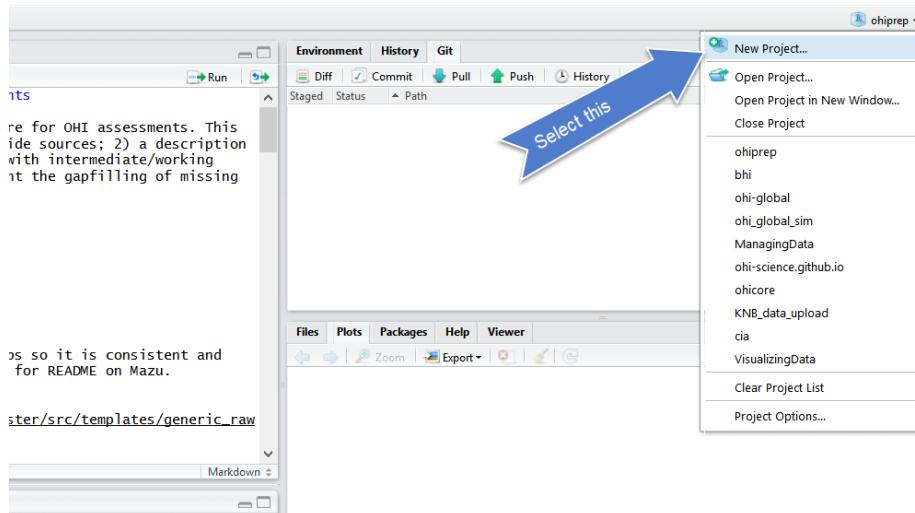
This will let us take advantage of something that is really key about GitHub.com: you can easily navigate through folders within repositories and the urls reflect this navigation. The greatness of this will be evident soon. So let's set ourselves up for easily translating (and remembering) those navigation paths by having a folder called `github` that will serve as our 'github.com'.

So really. Make sure that you have an all-lowercase folder called `github` in your home directory!!

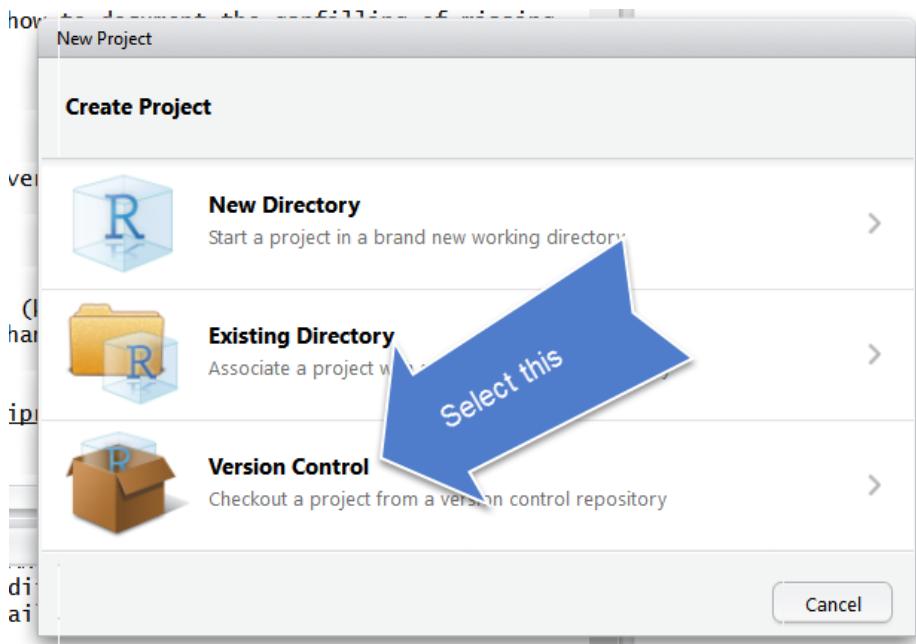
**Step 1:** Copy the web address of the repository you want to clone.



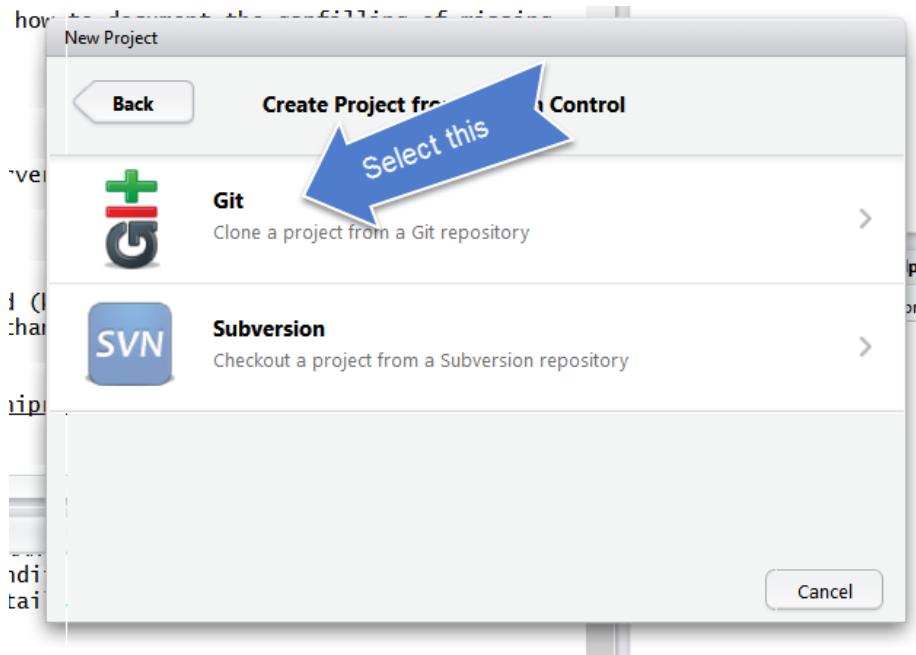
**Step 2:** from RStudio, go to New Project (also in the File menu).



**Step 3:** Select Version Control



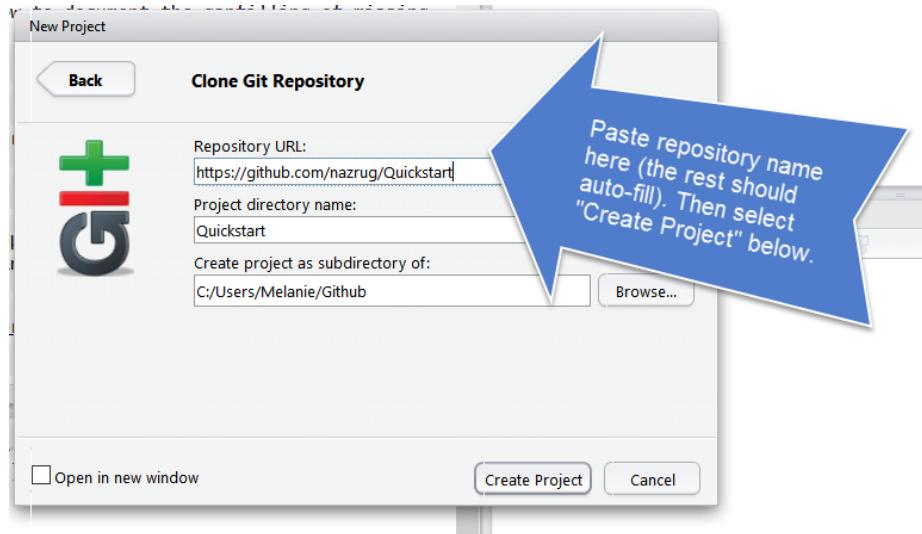
**Step 4:** Select Git



**Step 5:** Paste it in the Repository URL field, and type **tab** to autofill the Project Directory name. Make sure you keep the Project Directory Name **THE**

SAME as the repository name from the URL.

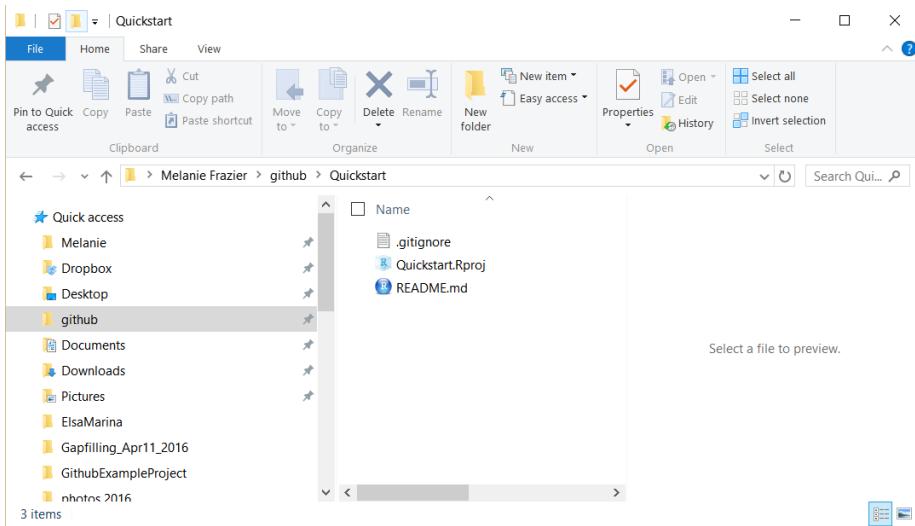
Save it in your github folder (click on Browse) to do this.



If everything went well, the repository will be added to the list located here:



And the repository will be saved to the Github folder on your computer:

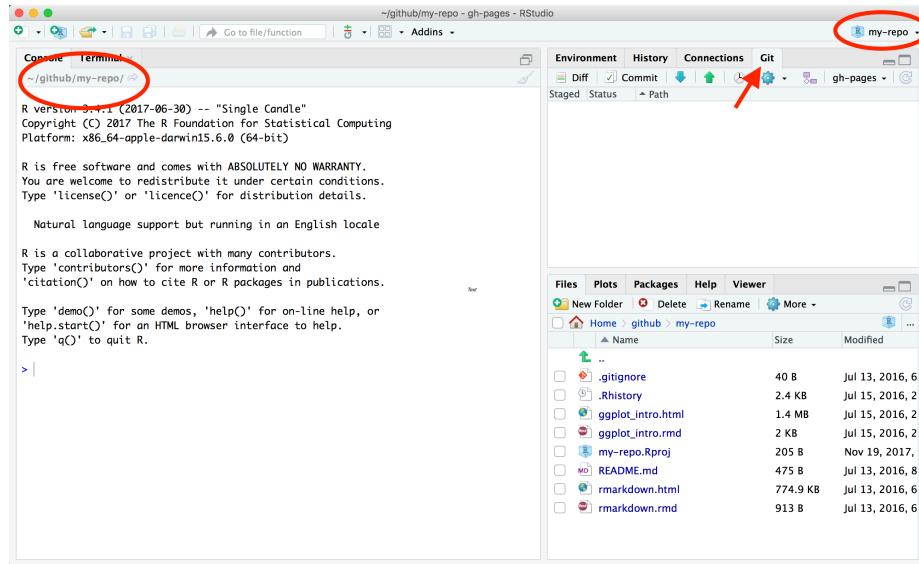


Ta da!!!! The folder doesn't contain much of interest, but we are going to change that.

## 4.8 Inspect your repository

Notice a few things in our repo here:

1. Our working directory is set to `~/github/r-workshop`. This means that I can start working with the files I have in here without setting the filepath. This is that when we cloned this from RStudio, it created an RStudio project, which you can tell because:
  - `.RProj` file, which you can see in the Files pane.
  - The project is named in the top right hand corner
2. We have a git tab! This is how we will interface directly to Github.com



When you first clone a repo through RStudio, RStudio will add an `.Rproj` file to your repo. And if you didn't add a `.gitignore` file when you originally created the repo on GitHub.com, RStudio will also add this for you. These will show up with little yellow ? icons in your git tab. This is GitHub's way of saying: "I am responsible for tracking everything that happens in this repo, but I haven't seen these files yet. Do you want me to track them too?" (We'll see that when you click the box to stage them, they will turn into As because they have been added to the repo.)

## 4.9 Add files to our local repo

The repository will contain:

- `.gitignore` file
- `README.md`
- `Rproj`

Let's create the following:

- folder called "data"
- folder called "figures"

They both show up in your Finder! ...

#### 4.9.1 Get data files into your working directory

In Session 1, we introduced how and why R Projects are great for reproducibility, because our self-contained working directory will be the **first** place R looks for files.

You downloaded several files for this workshop, some comma separate value (CSV) files and some as Excel spreadsheets (.xlsx):

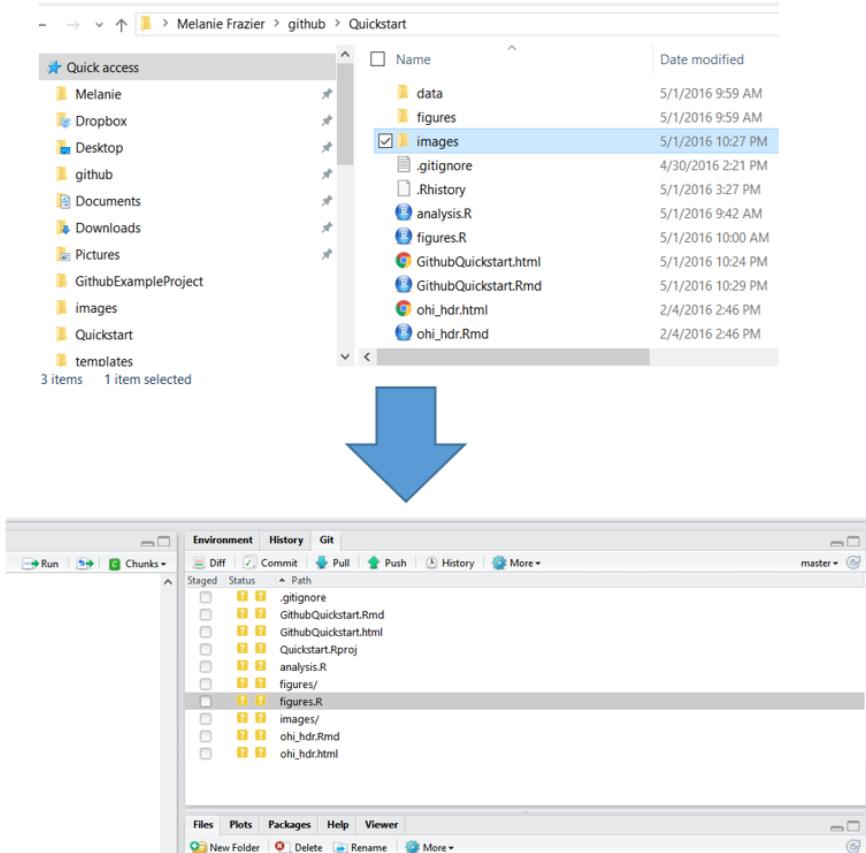
- fish\_counts\_curated.csv
- invert\_counts\_curated.xlsx
- kelp\_counts\_curated.xlsx
- substrate\_cover\_curated.xlsx
- lobsters.xlsx
- lobsters2.xlsx
- ca\_np.csv
- ci\_np.xlsx

Copy and paste those files into the ‘data’ subfolder of your R project. Notice that now these files are in your working directory when you go back to that Project in RStudio (check the ‘Files’ tab and navigate to the data subfolder). That means they’re going to be in the first place R will look when you ask it to find a file to read in.

I’m going to go to the Finder (Windows Explorer on a PC) and copy a file into my repository from there. And then I’m going to go back to RStudio – it shows up in the git tab! So the repository is being tracked, no matter how you make changes to it (changes do not have to be done only through RStudio).

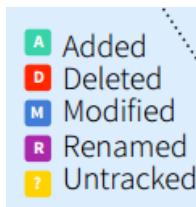
To make changes to the repository, you will work from your computer (“local Github”).

When files are changed in the local repository, these changes will be reflected in the Git tab of RStudio:



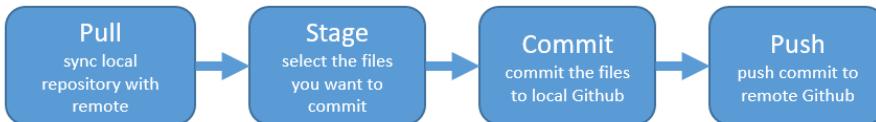
### 4.9.2 Inspect what has changed

These are the codes RStudio uses to describe how the files are changed, (from the RStudio cheatsheet):



## 4.10 Sync from RStudio to GitHub

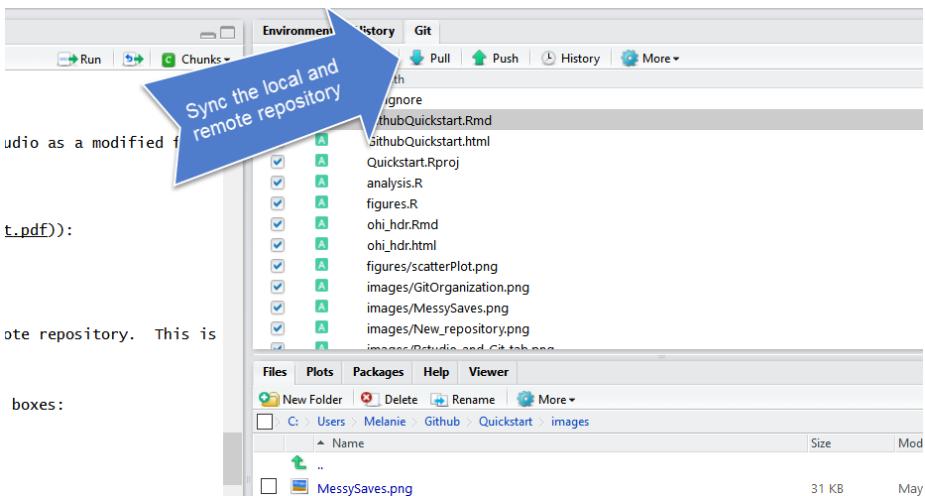
When you are ready to commit your changes, you follow these steps:



We walk through this process below:

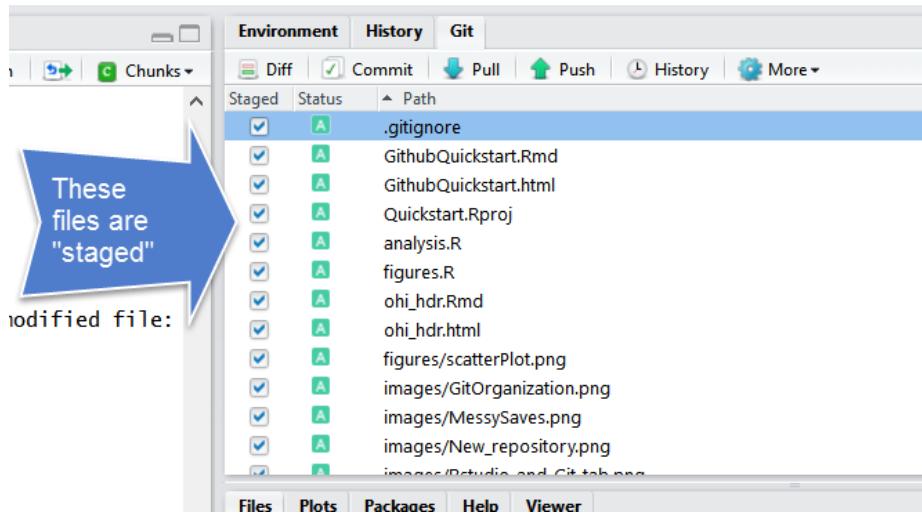
### 4.10.1 Pull

From the Git tab, “Pull” the repository. This makes sure your local repository is synced with the remote repository. This is very important if other people are making changes to the repository or if you are working from multiple computers.

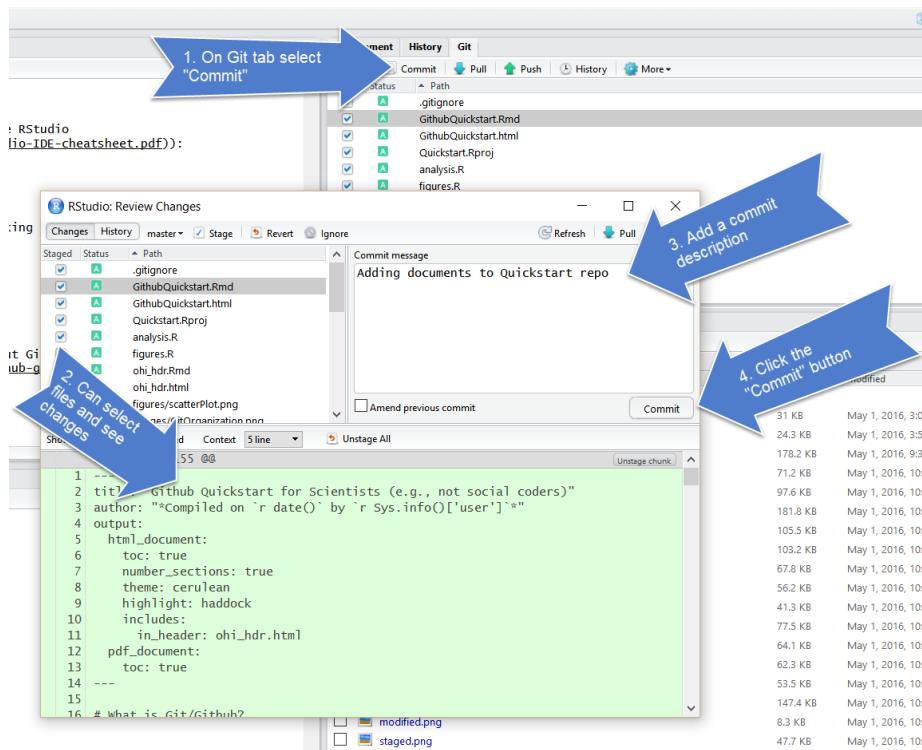


### 4.10.2 Stage

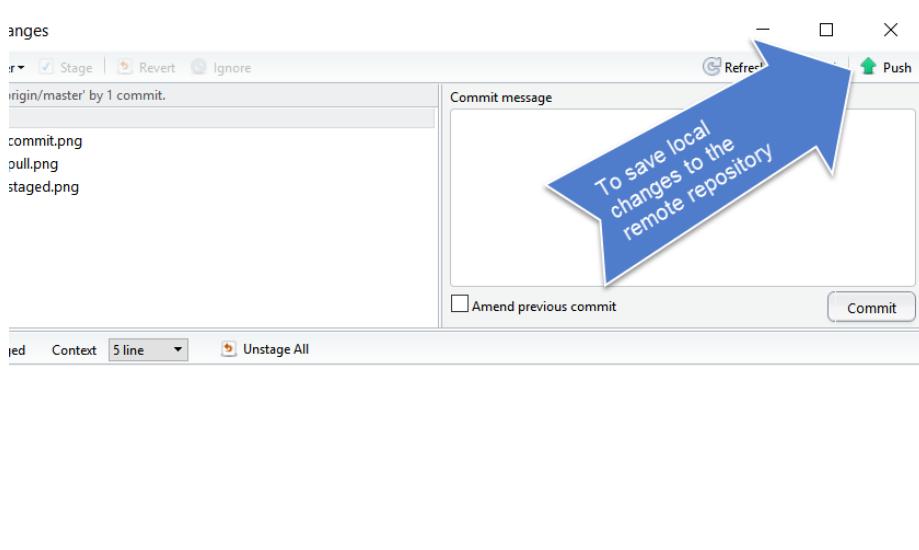
Stage the files you want to commit. In RStudio, this involves checking the “Staged” boxes:



### 4.10.3 Commit



#### 4.10.4 Push



### 4.11 Explore remote Github

The files you added should be on [github.com](https://github.com):

The screenshot shows a GitHub repository page for 'nazrug / Quickstart'. At the top, there are various browser extensions and tabs. Below the header, the repository name 'nazrug / Quickstart' is displayed along with a dropdown menu and several buttons: 'Unwatch', 'Star', 'Fork', and 'Settings'. The main content area shows the commit history:

File	Description	Time Ago
figures	Adding Quickstart files	21 minutes ago
images	Adding Quickstart files	21 minutes ago
.gitignore	Adding Quickstart files	21 minutes ago
GithubQuickstart.Rmd	Adding Quickstart files	21 minutes ago
GithubQuickstart.html	Adding Quickstart files	21 minutes ago
Quickstart.Rproj	Adding Quickstart files	21 minutes ago
README.md	Initial commit	2 days ago
analysis.R	Adding Quickstart files	21 minutes ago
figures.R	Adding Quickstart files	21 minutes ago
ohi_hdr.Rmd	Adding Quickstart files	21 minutes ago
ohi_hdr.html	Adding Quickstart files	21 minutes ago

Let's also explore commit history, file history.

#### 4.11.1 Activity

Go back to RStudio.

This time let's edit an existing file instead of adding something new. Open your README file by clicking on it in the Files pane (lower right corner). Write a few lines of text (like your dog's name), save, and see what happens in your Git Tab. Sync it to your remote repository at Github.com.

### 4.12 Create a new R Markdown file

Now get ourselves back into learning R. We are going to use R Markdown so that you can write notes to yourself in Markdown, and have a record of all your R code. Writing R commands in the console like we did this morning is great, but limited; it's hard to keep track of and hard to efficiently share with others. Plus, as your analyses get more complicated, you need to be able to see them all in one place.

Go to File > New File > R Markdown ... (or click the green plus in the top left corner).

Let's set up this file so we can use it for the rest of the day. I'm going to delete all the text that is already there and write some new text.

This will be your notes for the next session.

Here's what I'm going to write in my R Markdown file to begin:

```
---
```

```
title: "Reading data into R with `readxl`"
author: "Julie Lowndes"
date: "12/7/2019"
output: html_document
---
```

```
# Learning `readxl`
```

```
We are working with data and it's going to be amazing.
```

Now, let's save it. I'm going to call my file `readxl.Rmd`. You can knit it if you'd like.

Then, knit your file, and sync your file to GitHub: commit and pull

What if a file doesn't show up in the Git tab and you expect that it should? Check to make sure you've saved the file. If the filename is red with an asterix, there have been changes since it was saved. Remember to save before syncing to GitHub!

## 4.13 Explore your webpage

You've just created a webpage!

It will exist at this url: `username.github.io/repo-name/filename`. Mine is: `jules32.github.io/r-workshop/readxl`.

***ProTip*** Pay attention to URLs. An unsung skill of the modern analyst is to be able to navigate the internet by keeping an eye on patterns.

Troubleshooting:

- 404 error? Remove trailing / from the url
- Wants you to download? Remove trailing .Rmd from the url

## 4.14 Committing - how often? Tracking changes in your files

Whenever you make changes to the files in Github, you will walk through the Pull -> Stage -> Commit -> Push steps.

I tend to do this every time I finish a task (basically when I start getting nervous that I will lose my work). Once something is committed, it is very difficult to lose it.

One thing that I love about Github is that it is easy to see how files have changed over time. Usually I compare commits through [github.com](https://github.com):

The screenshot shows two views of a GitHub repository's commit history. The top view shows a list of 1,662 commits across 2 branches and 1 release, with 8 contributors. A blue arrow points to the first commit, which has a message: "debugged and reprocessed the SPP goal... done with that? now to make ...". The bottom view is a zoomed-in look at the first commit, showing the full message: "debugged and reprocessed the SPP goal... done with that? now to make ... committed 7 hours ago". A large blue arrow points to this commit message. Another blue arrow points to the short commit ID "5f8aba9" located below the message.

Date	Message	Time Ago
May 2, 2016	debugged and reprocessed the SPP goal... done with that? now to make ...	committed 7 hours ago
May 1, 2016	Merge branch 'master' of https://github.com/OHI-Science/ohiprep debugging the species for 2016... I think it's actually OK...	committed a day ago
Apr 28, 2016	Create README.md data organization SOP revisions	committed 4 days ago
Apr 27, 2016		

You can click on the commits to see how the files changed from the previous commit:

The screenshot shows a GitHub pull request interface. The title of the pull request is "debugging the species for 2016... I think it's actually OK...". The commit message is "= committed a day ago". The commit hash is b21fe2946589b4b2ed7351bc60c154651aab7953. The commit message indicates 2,432,873 additions and 426,265 deletions. A note states "Sorry, we could not display the entire diff because it was too big." Below this, the diff view shows a file named "globalprep/spp\_ico/v2016/data\_prep\_SPP.Rmd". The diff highlights deleted code in red and added code in green. Two large blue arrows point from the text "Code in red was deleted" and "Code in green was added" towards the respective colored regions in the diff view.

```

@@ -22,18 +22,18 @@ library(Foreign)
22 22   library(data.table)
23 23   library(sp)
24 24   library(rgdal)
25 25 -library(raster)
26 26 +library(raster)
27 27 +library(raster)
28 28
29 29   source('~/github/ohiprep/src/R/common.R')
30 30
31 31   goal <- 'globalprep/spp_ico'
32 32   scenario <- 'v2016'
33 33   dir_anx <- file.path(dir_M, 'git-annex', goal)
34 34 -dir_data_am <- file.path(dir_M, 'git-annex/globalprep/_raw_data', 'aquamaps', str_replace(scenario, 'v', 'd'))
35 35 -dir_data_iucn <- file.path(dir_M, 'git-annex/globalprep/_raw_data', 'iucn_spp')
36 36 -dir_data_bird <- file.path(dir_M, 'git-annex/globalprep/_raw_data', 'birdlife_intl')
37 37 +dir_data_am <- file.path(dir_M, 'git-annex/globalprep/_raw_data', 'aquamaps', d2015)
38 38 +dir_data_iucn <- file.path(dir_M, 'git-annex/globalprep/_raw_data', 'iucn_spp', d2015)
39 39 +dir_data_bird <- file.path(dir_M, 'git-annex/globalprep/_raw_data', 'birdlife_intl', d2015)
40 40
41 41   dir_git <- file.path('~/github/ohiprep', goal)
42 42
43 43   #remove old annex data files scenario from RDS file

```

## 4.15 Happy Git with R

If you have problems, we'll help you out using Jenny Bryan's HappyGitWithR, particularly the sections on Detect Git from RStudio and RStudio, Git, GitHub Hell (troubleshooting). So as we are coming around, have a look at it and see if you can help troubleshoot too!

## 4.16 Efficiency Tips



# Chapter 5

## Graphs with ggplot2

### 5.1 Summary

Now that we know how to *get* some data, the next thing we'll probably want to do is look at it. In Excel, graphs are made by manually selecting options - which, as we've discussed previously, may not be the best option for reproducibility. Also, if we haven't built a graph with reproducible code, then we might not be able to easily recreate a graph *or* use that code again to make the same style graph with different data.

Using `ggplot2`, the graphics package within the `tidyverse`, we'll write reproducible code to manually and thoughtfully build our graphs.

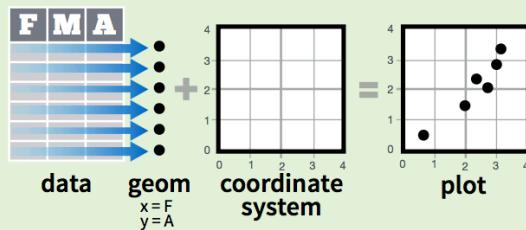
“`ggplot2` implements the grammar of graphics, a coherent system for describing and building graphs. With `ggplot2`, you can do more faster by learning one system and applying it in many places.” - R4DS

So yeah...that `gg` is from “grammar of graphics.”

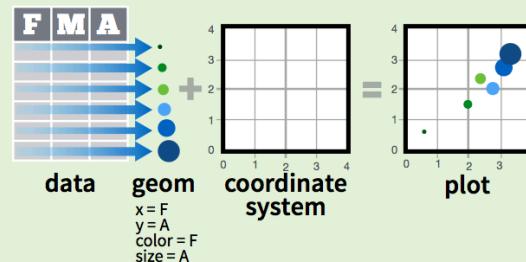
We'll use the `ggplot2` package, but the function we use to initialize a graph will be `ggplot`, which works best for data in tidy format (i.e., a column for every variable, and a row for every observation). Graphics with `ggplot` are built step-by-step, adding new elements as layers with a plus sign (+) between layers (note: this is different from the pipe operator, `%>%`. Adding layers in this fashion allows for extensive flexibility and customization of plots.

## Basics

**ggplot2** is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data** set, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.



To display data values, map variables in the data set to aesthetic properties of the geom like **size**, **color**, and **x** and **y** locations.



### 5.1.1 Objectives

- Build several common types of graphs (scatterplot, column, line) in ggplot2
- Customize gg-graph aesthetics (color, style, themes, etc.)
- Update axis labels and titles
- Combine compatible graph types (geoms)
- Build multiseries graphs
- Split up data into faceted graphs
- Exporting figures with `ggsave()`

### 5.1.2 Resources

- <https://r4ds.had.co.nz/data-visualisation.html>
- ggplot2-cheatsheet-2.1.pdf
- Graphs with ggplot2 - Cookbook for R
- “Why I use ggplot2” - David Robinson Blog Post

## 5.2 Getting started - Create a new .Rmd, attach packages & get data

Within your existing version-controlled R project, create a new R Markdown document with title “Data visualization with ggplot2.” Remove everything below the first code chunk. Knit and save the .Rmd file within your project working directory as “my\_ggplot2”.

The `ggplot2` package is part of the `tidyverse`, so we don’t need to attach it separately. Attach the `tidyverse`, `readxl` and `here` packages in the top-most code chunk of your .Rmd.

```
library(tidyverse)
library(readxl)
library(here)
```

In this session, we’ll use data for parks visitation from two files:

- A comma-separated-value (CSV) file containing visitation data for all National Parks in California (`ca_np.csv`)
- A single Excel worksheet containing only visitation for Channel Islands National Park (`ci_np.xlsx`)

Add a new code chunk to read in the data from the `data` subfolder within your working directory.

```
ca_np <- read_csv(here("data", "ca_np.csv"))
ci_np <- read_xlsx(here("data", "ci_np.xlsx"))
```

Let’s take a quick look at the data to see what it contains. For example:

- `View()`: to look at the object in spreadsheet format
- `names()`: to see the variable (column) names
- `summary()`: see a quick summary of each variable

### 5.3 Our first ggplot graph: Visitors to Channel Islands NP

To create a bare-bones ggplot graph, we need to tell R three basic things:

1. We're using `ggplot`
2. Data we're using & variables we're plotting (i.e., what is x and/or y?)
3. What type of graph we're making (the type of `geom`)

Generally, that structure will look like this:

```
ggplot(data = df_name, aes(x = x_var_name, y = y_var_name)) +
  geom_type()
```

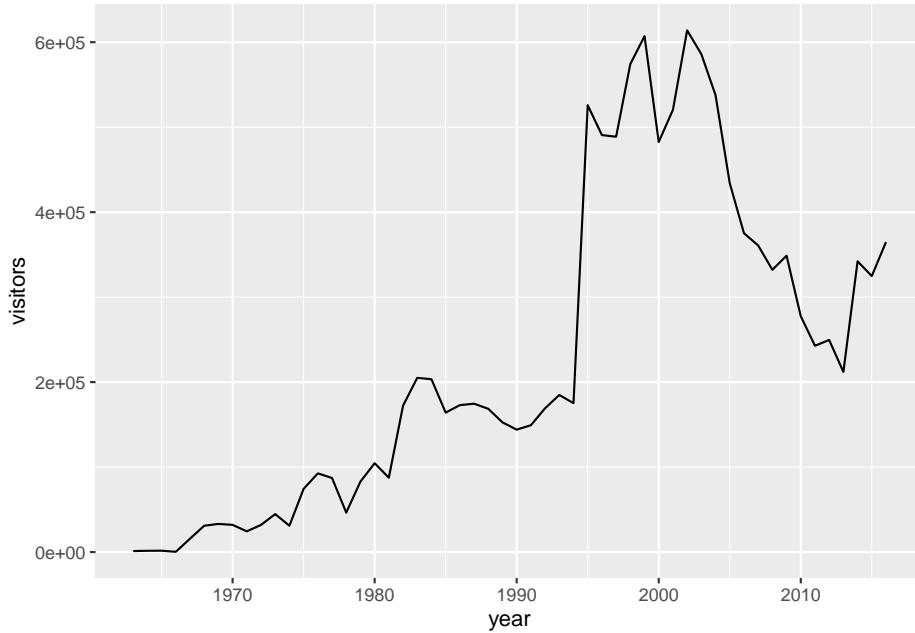
Breaking that down:

- First, tell R you're using `ggplot()`
- Then, tell it the object name where variables exist (`data = df_name`)
- Next, tell it the aesthetics `aes()` to specify which variables you want to plot
- Then add a layer for the type of geom (graph type) with `geom_*`() - for example, `geom_point()` is a scatterplot, `geom_line()` is a line graph, `geom_col()` is a column graph, etc.

Let's do that to create a line graph of visitors to Channel Islands National Park:

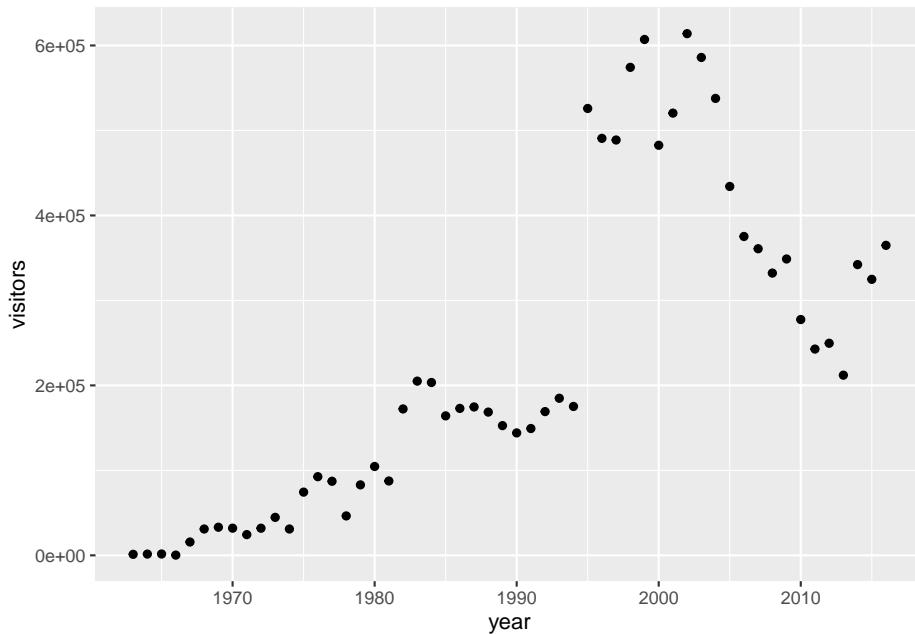
```
ggplot(data = ci_np, aes(x = year, y = visitors)) +
  geom_line()
```

### 5.3. OUR FIRST GG PLOT GRAPH: VISITORS TO CHANNEL ISLANDS NP57



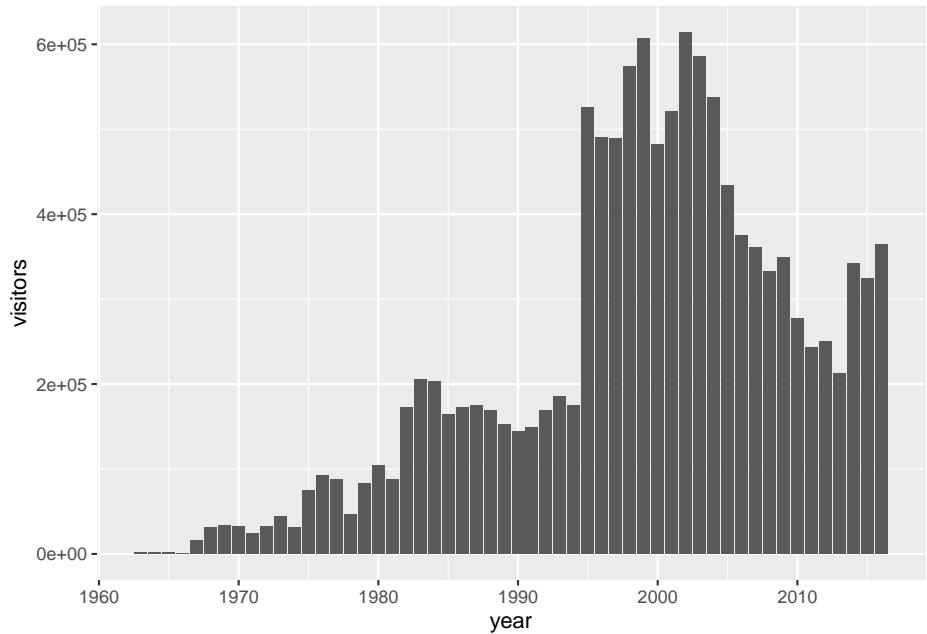
Or, we could change that to a scatterplot just by updating the `geom_*`:

```
ggplot(data = ci_np, aes(x = year, y = visitors)) +  
  geom_point()
```



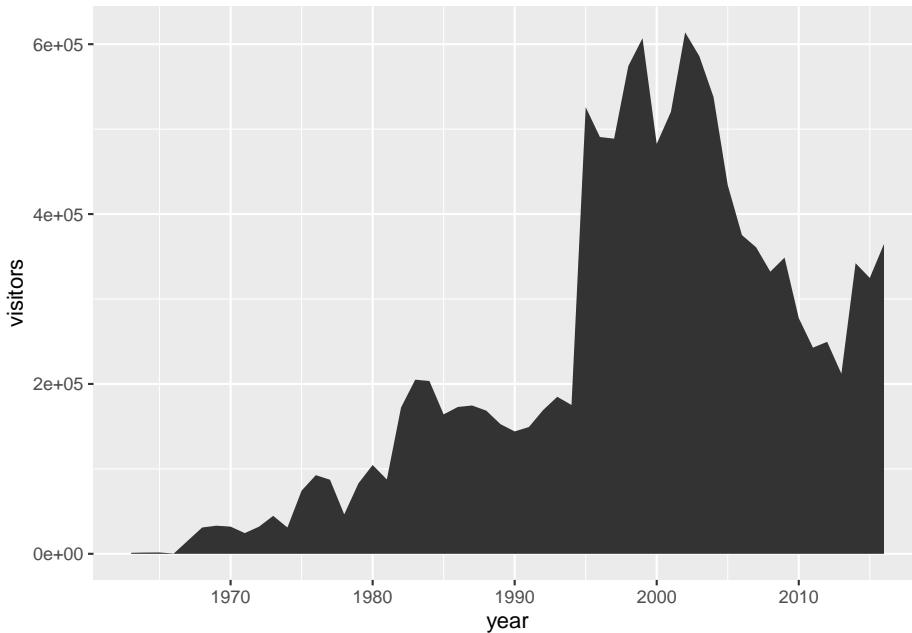
We could even do that for a column graph:

```
ggplot(data = ci_np, aes(x = year, y = visitors)) +  
  geom_col()
```



Or an area plot...

```
ggplot(data = ci_np, aes(x = year, y = visitors)) +  
  geom_area()
```



We can see that updating to different `geom_*` types is quick, so long as the types of graphs we're switching between are compatible.

The data are there, now let's do some data viz customization.

## 5.4 Intro to customizing ggplot graphs

First, we'll customize some aesthetics (e.g. colors, styles, axis labels, etc.) of our graphs based on non-variable values.

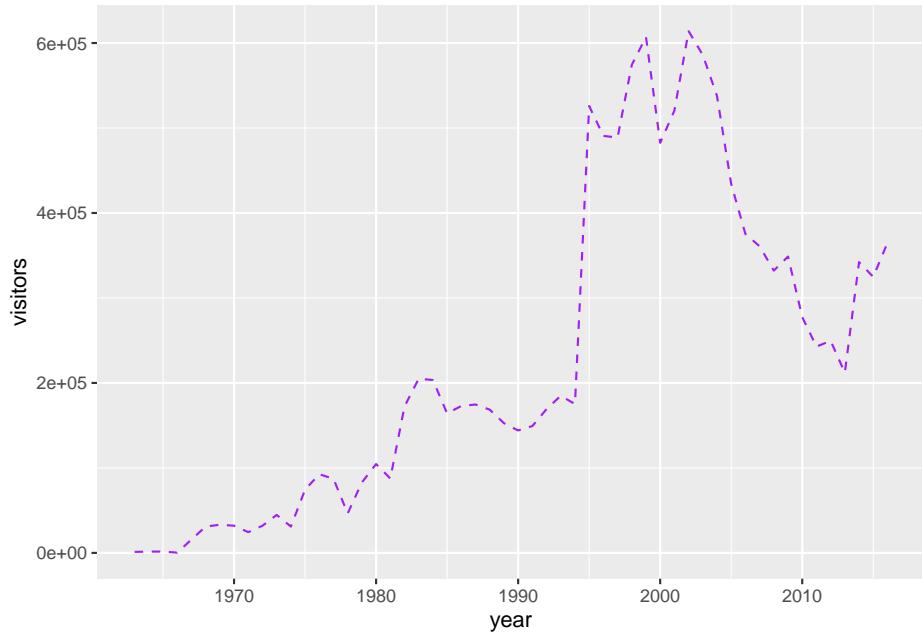
We can change the aesthetics of elements in a ggplot graph by adding arguments within the layer where that element is created.

Some common arguments we'll use first are:

- `color` = or `colour` =: update point or line colors
- `fill` =: update fill color for objects with areas
- `linetype` =: update the line type (dashed, long dash, etc.)
- `pch` =: update the point style
- `size` =: update the element size (e.g. of points or line thickness)
- `alpha` =: update element opacity (1 = opaque, 0 = transparent)

Building on our first line graph, let's update the line color to “purple” and make the line type “dashed”:

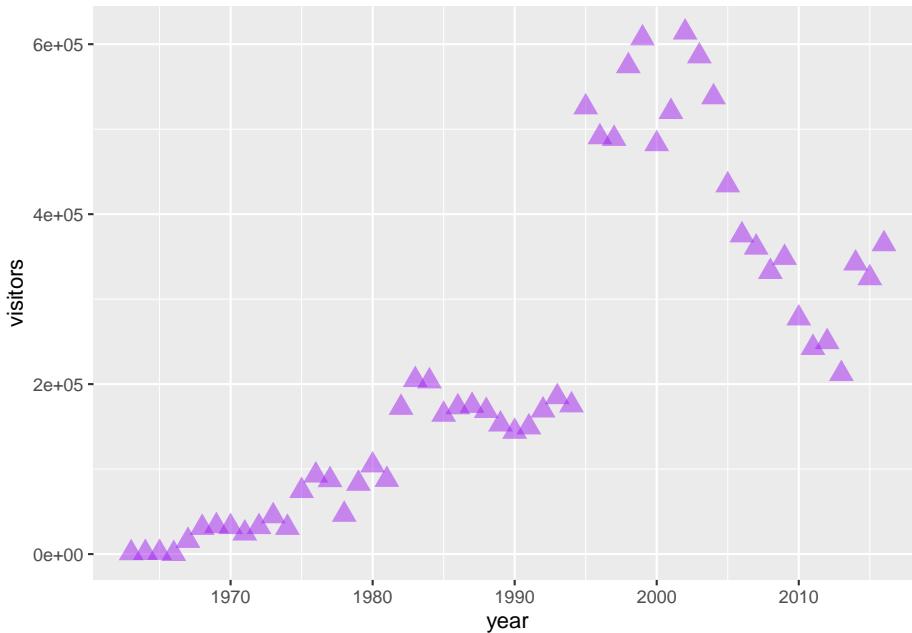
```
ggplot(data = ci_np, aes(x = year, y = visitors)) +
  geom_line(
    color = "purple",
    linetype = "dashed"
  )
```



How do we know which color names ggplot will recognize? If you google “R colors ggplot2” you’ll find a lot of good resources. Here’s one: SAPE ggplot2 colors quick reference guide

Now let’s update the point, style and size of points on our previous scatterplot graph using `color =`, `size =`, and `pch =` (see `?pch` for the different point styles, which can be further customized).

```
ggplot(data = ci_np, aes(x = year, y = visitors)) +
  geom_point(color = "purple",
             pch = 17,
             size = 4,
             alpha = 0.5)
```



#### 5.4.1 Activity: customize your own ggplot graph

Update one of the example graphs you created above to customize **at least** an element color and size!

## 5.5 Mapping variables onto aesthetics

In the examples above, we have customized aesthetics based on constants that we input as arguments (e.g., the color / style / size isn't changing based on a variable characteristic or value). Sometimes, however, we **do** want the aesthetics of a graph to depend on a variable. To do that, we'll **map variables onto graph aesthetics**, meaning we'll change how an element on the graph looks based on a variable characteristic (usually, character or value).

When we want to customize a graph element based on a variable's characteristic or value, add the argument within `aes()` in the appropriate `geom_*`() layer

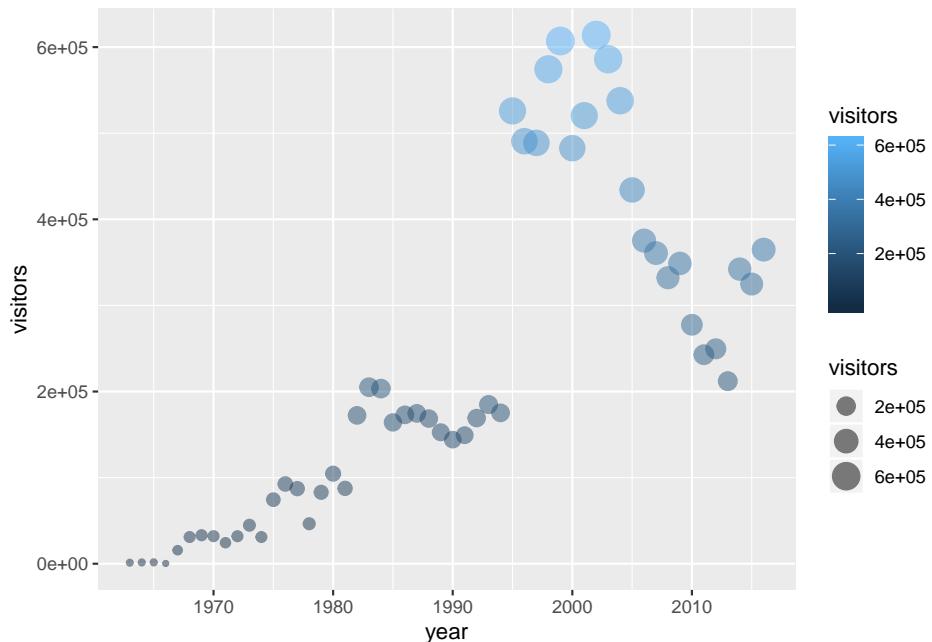
In short, if updating aesthetics based on a variable, make sure to put that argument inside `aes()`.

**Example:** Create a ggplot scatterplot graph where the **size** and **color** of the points change based on the **number of visitors**, and make all points the same

level of opacity (`alpha = 0.5`). Notice the `aes()` around the `size =` and `color =` arguments.

Also: this is overmapped and unnecessary. Avoid excessive / overcomplicated aesthetic mapping in data visualization.

```
ggplot(data = ci_np, aes(x = year, y = visitors)) +
  geom_point(
    aes(size = visitors,
        color = visitors),
    alpha = 0.5
  )
```

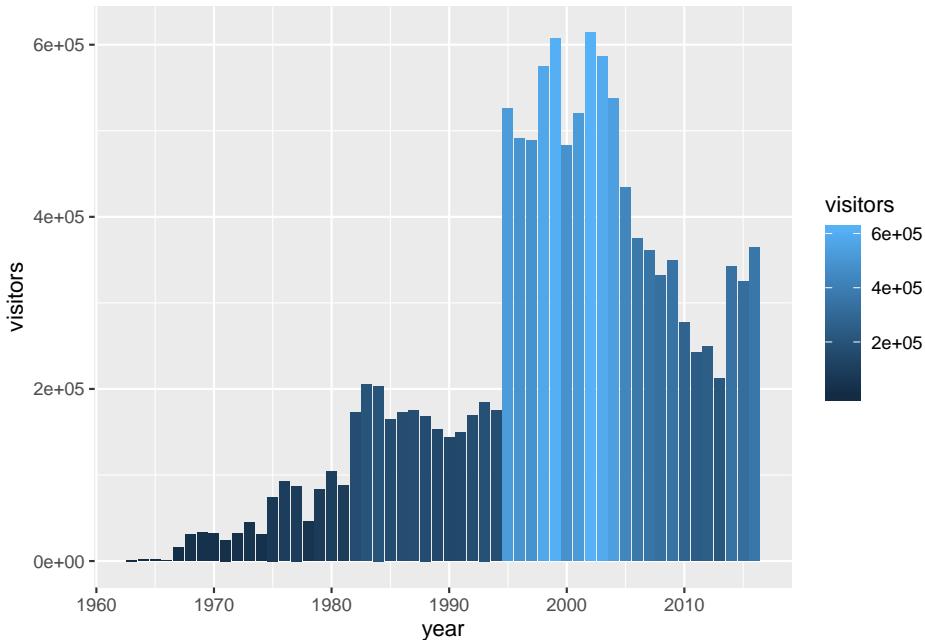


In the example above, notice that the two arguments that `do` depend on variables are within `aes()`, but since `alpha = 0.5` doesn't depend on a variable then it is *outside the `aes()` but still within the `geom_point()` layer*.

### 5.5.1 Activity: map variables onto graph aesthetics

Create a column plot of Channel Islands National Park visitation over time, where the **fill color** (argument: `fill =`) changes based on the number of **visitors**.

```
ggplot(data = ci_np, aes(x = year, y = visitors)) +
  geom_col(aes(fill = visitors))
```



Sync your project with your GitHub repo.

## 5.6 ggplot2 complete themes

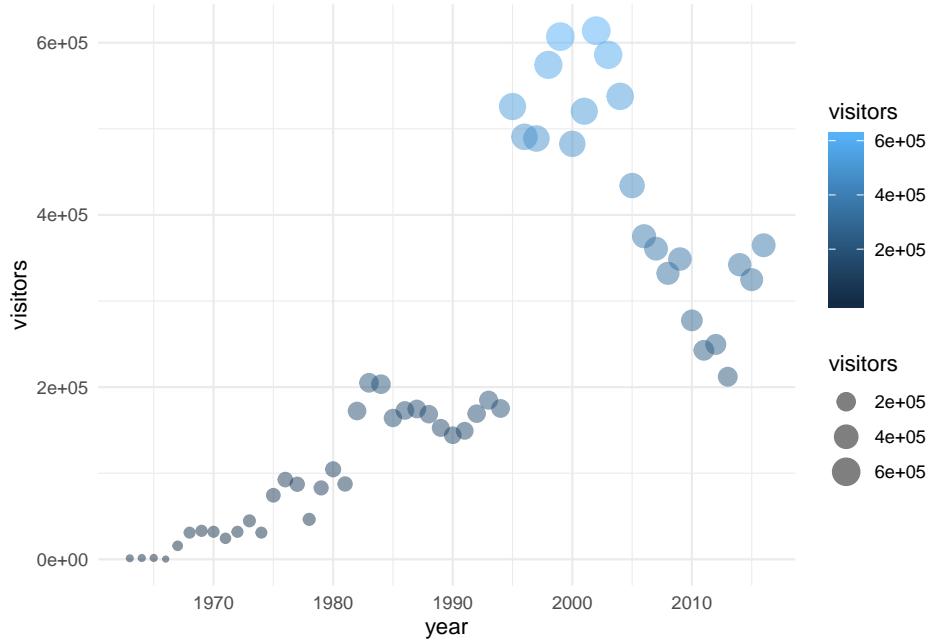
While every element of a ggplot graph is manually customizable, there are also built-in themes (`theme_*`) that you can add to your ggplot code to make some major headway before making smaller tweaks manually.

Here are a few to try today (but also notice all the options that appear as we start typing `theme_` into our ggplot graph code!):

- `theme_light()`
- `theme_minimal()`
- `theme_bw()`

Here, let's update our previous graph with `theme_minimal()`:

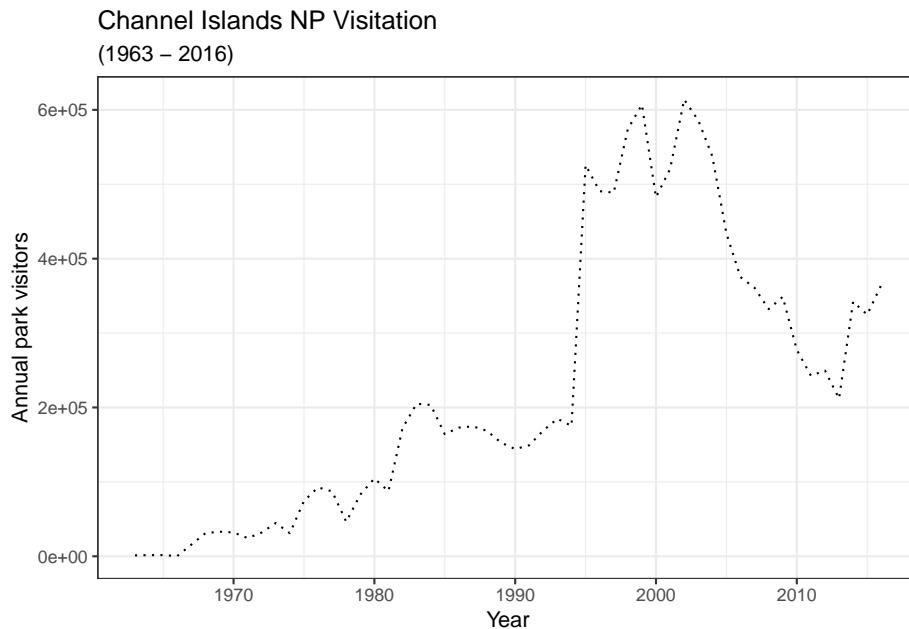
```
ggplot(data = ci_np, aes(x = year, y = visitors)) +
  geom_point(
    aes(size = visitors,
        color = visitors),
    alpha = 0.5
  ) +
  theme_minimal()
```



## 5.7 Updating axis labels and titles

Use `labs()` to update axis labels, and add a title and/or subtitle to your ggplot graph.

```
ggplot(data = ci_np, aes(x = year, y = visitors)) +
  geom_line(linetype = "dotted") +
  theme_bw() +
  labs(
    x = "Year",
    y = "Annual park visitors",
    title = "Channel Islands NP Visitation",
    subtitle = "(1963 - 2016)"
  )
```



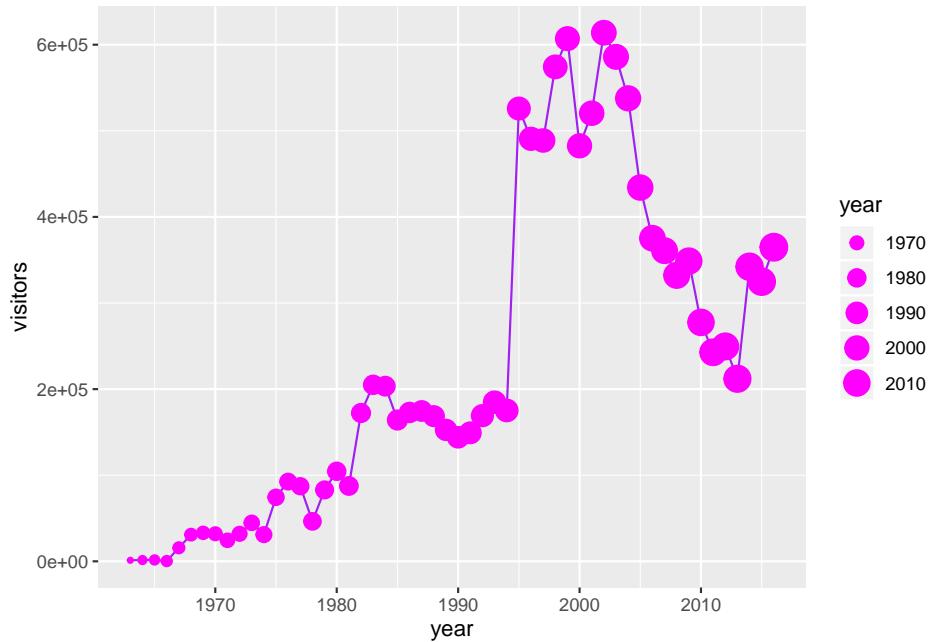
**Note:** If you want to update the formatting of axis values (for example, to convert to comma format instead of scientific format above), you can use the `scales` package options (see more from the R Cookbook).

## 5.8 Combining compatible geoms

As long as the geoms are compatible, we can layer them on top of one another to further customize a graph.

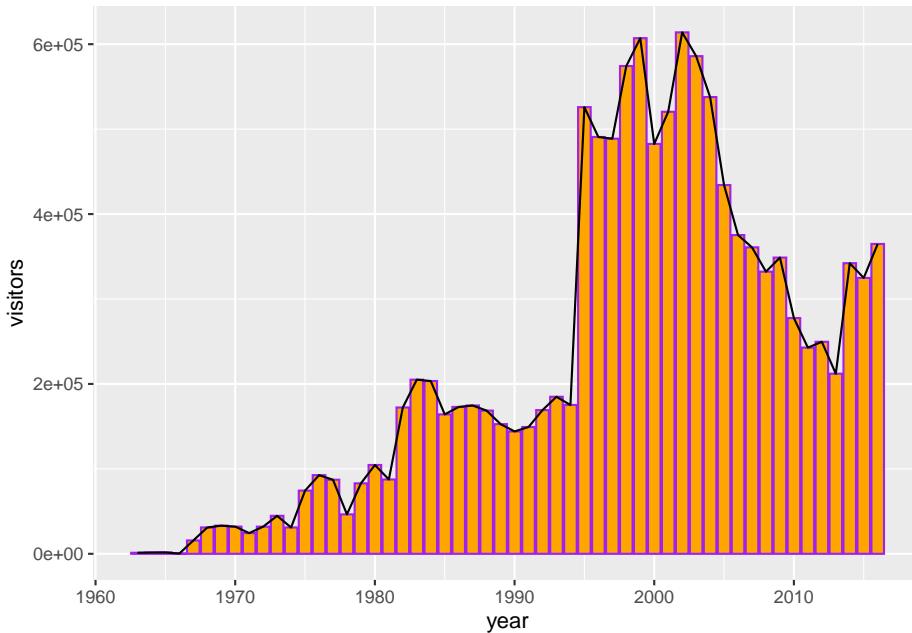
For example, adding points to a line graph:

```
ggplot(data = ci_np, aes(x = year, y = visitors)) +
  geom_line(color = "purple") +
  geom_point(color = "magenta",
             aes(size = year))
```



Or, combine a column and line graph (not sure why you'd want to do this, but you can):

```
ggplot(data = ci_np, aes(x = year, y = visitors)) +
  geom_col(fill = "orange",
           color = "purple") +
  geom_line()
```

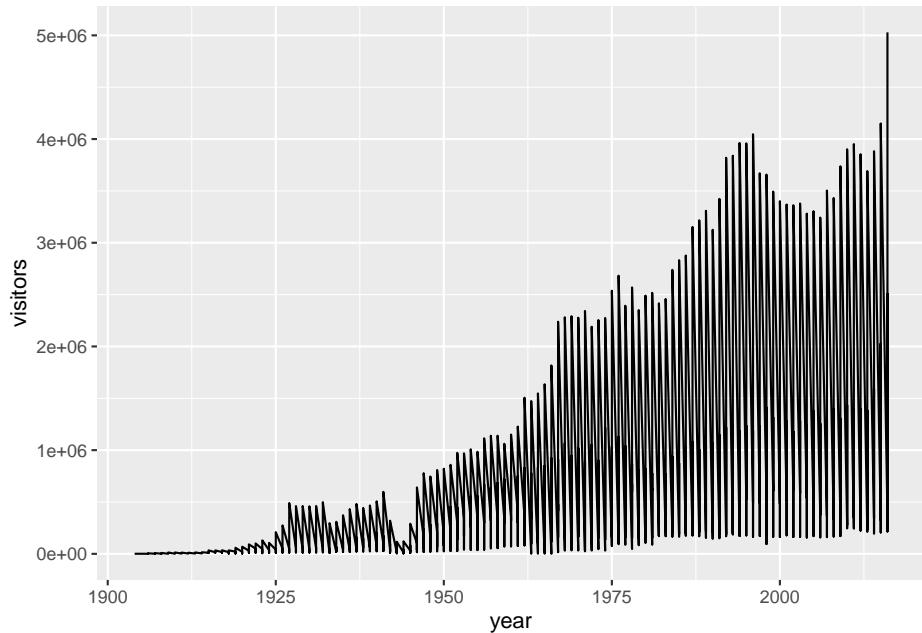


## 5.9 Multi-series ggplot graphs

In the examples above, we only had a single series - visitation at Channel Islands National Park. Often we'll want to visualize multiple series. For example, from the `ca_np` object we have stored, we might want to plot visitation for *all* California National Parks.

To do that, we need to add an aesthetic that lets `ggplot` know how things are going to be grouped. A demonstration of why that's important - what happens if we *don't* let `ggplot` know how to group things?

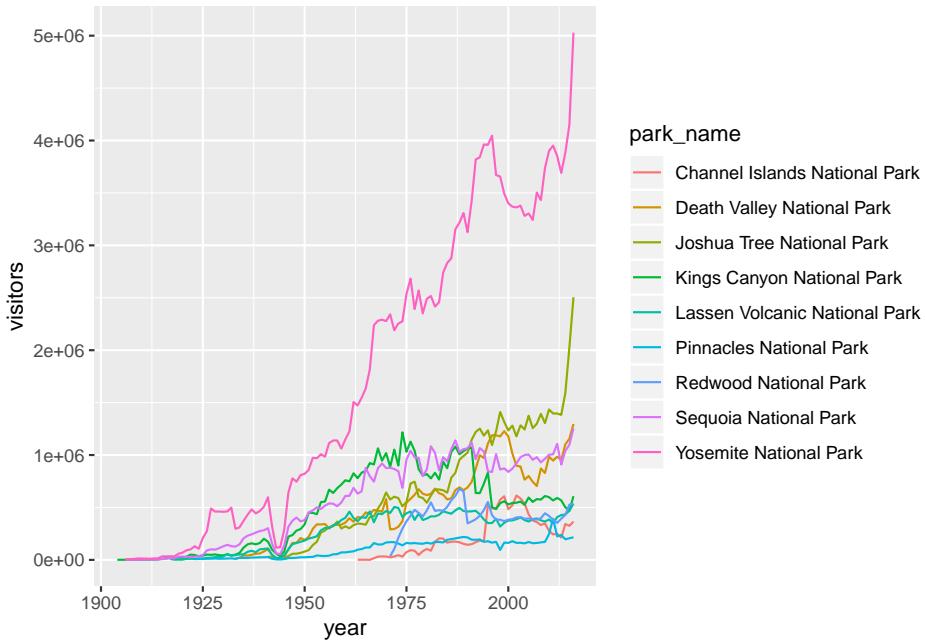
```
ggplot(data = ca_np, aes(x = year, y = visitors)) +
  geom_line()
```



Well that's definitely a mess, and it's because ggplot has no idea that these **should be different series based on the different parks that appear in the 'park\_name' column.**

We can make sure R does know by updating an aesthetic based on *park\_name*:

```
ggplot(data = ca_np, aes(x = year, y = visitors, color = park_name)) +  
  geom_line()
```



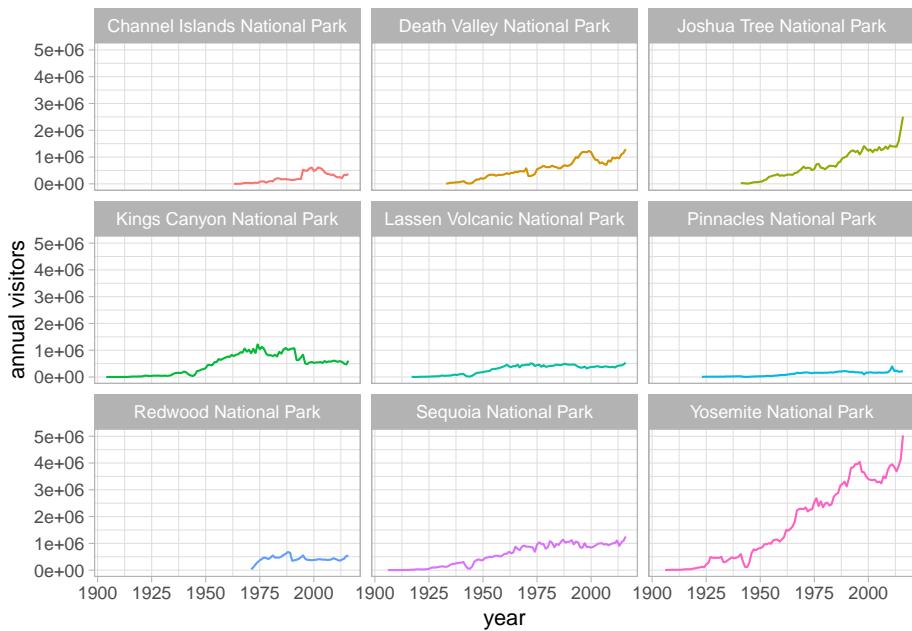
**Note:** You could also add that aesthetic (`color = park_name`) in the `geom_line()` layer, instead of in the topmost `ggplot()` layer.

## 5.10 Faceting ggplot graphs

When we facet graphs, we split them up into multiple plotting panels, where each panel contains a subset of the data. In our case, we'll split the graph above into different panels, each containing visitation data for a single park.

Also notice that any general theme changes made will be applied to *all* of the graphs.

```
ggplot(data = ca_np, aes(x = year, y = visitors, color = park_name)) +
  geom_line(show.legend = FALSE) +
  theme_light() +
  labs(x = "year", y = "annual visitors") +
  facet_wrap(~ park_name)
```



## 5.11 Exporting a ggplot graph with ggsave()

If we want our graph to appear in a knitted html, then we don't need to do anything else. But often we'll need a saved image file, of specific size and resolution, to share or for publication.

`ggsave()` will export the *most recently run* ggplot graph by default (`plot = last_plot()`), unless you give it the name of a different saved ggplot object. Some common arguments for `ggsave()`:

- `width` := set exported image width (default inches)
- `height` := set exported image height (default height)
- `dpi` := set dpi (dots per inch)

So to export the faceted graph above at 180 dpi, width a width of 8" and a height of 7", we can use:

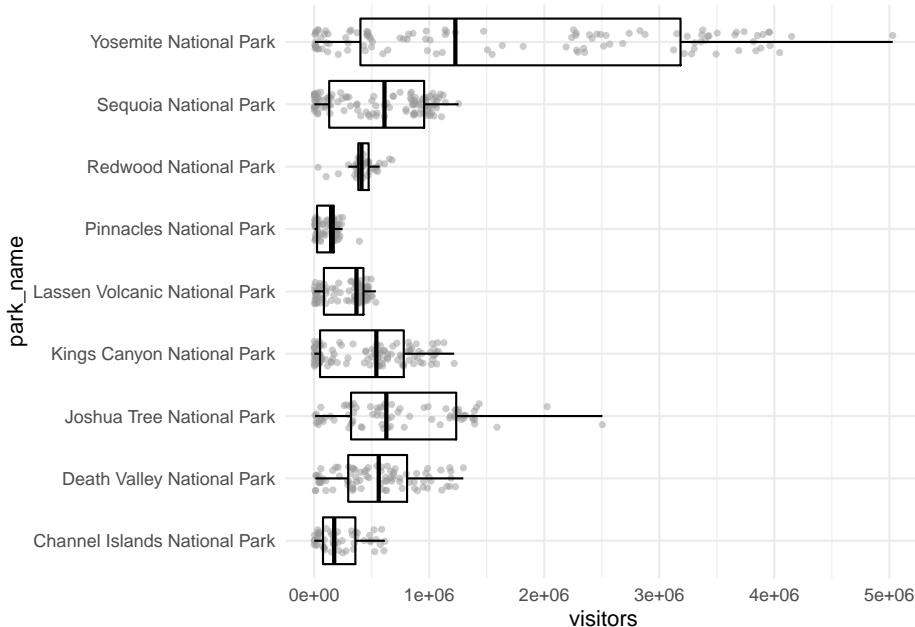
```
ggsave(here("figures", "np_graph.jpg"), dpi = 180, width = 8, height = 7)
```

Notice that a .jpg image of that name and size is now stored in your project working directory. You can change the type of exported image, too (e.g. pdf, tiff, eps, png, mmp, svg).

### 5.11.1 One final graph example: jitter and boxplots

For the record: this is not a good option for showing the visitation data because values are not independent observations of a random variable. But, for the purposes of showing a graph, we'll use visitation as our continuous measured variable in a jitter + boxplot anyway.

```
ggplot(data = ca_np, aes(x = park_name, y = visitors)) +
  geom_jitter(alpha = 0.5,
              color = "gray60",
              width = 0.2,
              size = 1) +
  geom_boxplot(fill = NA,
               color = "black",
               outlier.color = NA) +
  coord_flip() +
  theme_minimal()
```



Sync your project with your GitHub repo.

## 5.12 End ggplot session



# Chapter 6

## dplyr and Pivot Tables

### 6.1 Summary

Pivot tables are powerful tools in Excel for summarizing data in different ways. We will create these tables using the `group_by` and `summarize` functions from the `dplyr` package (part of the Tidyverse). We will also learn how to format tables and practice creating a reproducible report using RMarkdown and sharing it with GitHub.

### 6.2 Objectives

In R, we can use `dplyr` for pivot tables by using 2 main verbs in combination: `group_by` and `summarize`. We will also continue to emphasize reproducibility in all our analyses.

- Discuss pivot tables in Excel
- Introduce `group_by()` `%>%` `summarize()` from the `dplyr` package
- Practice our reproducible workflow with RMarkdown and GitHub

### 6.3 Resources

- [dplyr.tidyverse.org](http://dplyr.tidyverse.org)
- R for Data Science: Transform Chapter
- Intro to Pivot Tables I-III by Excel Campus (YouTube)

## 6.4 Pivot table overview

Wikipedia describes a pivot table as a “table of statistics that summarizes the data of a more extensive table...This summary might include sums, averages, or other statistics, which the pivot table groups together in a meaningful way.” Fun fact: it also says that “Although pivot table is a generic term, Microsoft trademarked PivotTable in the United States in 1994.”

Pivot tables are a really powerful tool for summarizing data, and we can have similar functionality in R — as well as nicely automating and reporting these tables. We will learn about this using data about lobsters and will go back and forth between R and Excel as we learn.

Let’s start off in R, and have a look at the data.

## 6.5 RMarkdown setup

Let’s start a new RMarkdown file in our repo, at the top-level (where it will be created by default in our Project). I’ll call mine `pivot_lobsters.Rmd`.

In the setup chunk, let’s attach our libraries and read in our lobster data. In addition to the `tidyverse` package we will also use the `skimr` package. You will have to install it, but don’t want it to be installed every time you write your code. The following is a nice convention for having the install instructions available (on the same line) as the `library()` call.

```
## attach libraries
library(tidyverse)
library(readxl)
library(here)
library(skimr) # install.packages('skimr')

## read in data
lobsters <- read_xlsx(here("data/lobsters.xlsx"))
```

Let’s add a code chunk and explore the data in a few ways.

```
# explore data
head(lobsters) # year and month as well as a column for date

## # A tibble: 6 x 7
##   year month date    site transect replicate size_mm
##   <dbl> <dbl> <chr>  <chr>   <dbl> <chr>      <dbl>
## 1  2012     8 8/20/12 ivee       3 A           70
## 2  2012     8 8/20/12 ivee       3 B           60
## 3  2012     8 8/20/12 ivee       3 B           65
## 4  2012     8 8/20/12 ivee       3 B           70
```

```
## 5 2012     8 8/20/12 ivee      3 B      85
## 6 2012     8 8/20/12 ivee      3 C      60
```

`head()` gives us a look at the first rows of the data (6 by default). I like this because I can see the column names and get a sense of the shape of the data. I can also see the class of each column (double or character)

In this data set, every row is a unique observation. This is called “uncounted” data; you’ll see there is no row for how many lobsters were seen because each row is an observation, or an “n of 1”.

```
# explore data
summary(lobsters)
```

`summary` gives us summary statistics for each variable (column). I like this for numeric columns, but it doesn’t give a lot of useful information for non-numeric data. To have a look there I like using the `skimr` package:

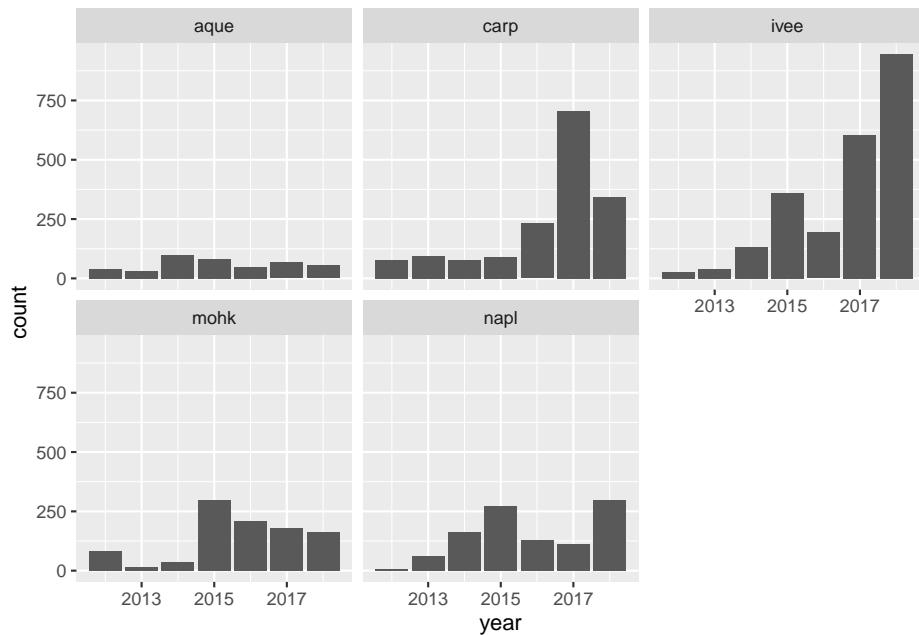
```
# explore data
skimr::skim(lobsters)
```

This `skimr::` notation is a reminder to me that `skim` is from the `skimr` package. It is a nice convention: it’s a reminder to others (especially you!).

`skim` lets us look more at each variable. I particularly like looking at missing data. There are 6 missing values in the `size_mm` variable.

We can also make a quick plot to have a look at these data, and use our new `ggplot2` skills. Let’s make a bar chart by year for each site

```
ggplot(lobsters, aes(x = year)) +
  geom_bar() +
  facet_wrap(~site)
```



(geom\_bar() counts things and geom\_col() is for values within the data (mean))

### 6.5.1 Our task

So this is all great to get a quick look. But what if we needed to report to someone about how the average size of lobsters has changed over time across sites?

To answer this we need to do a pivot table in Excel, or data wrangling in R.

Let's start by having a quick look at what pivot tables can do in Excel.

## 6.6 Pivot table demo

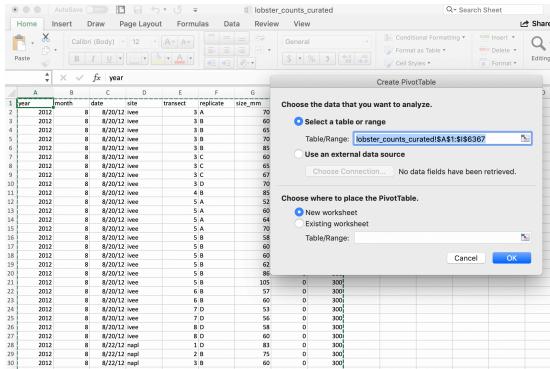
Let's make a pivot table with our lobster data.

Let's start off with how many lobsters were counted each year. I want a count of rows by year.

So to do this in Excel we would initiate the Pivot Table Process:

year	month	date	site	transect	replicate
2012	8	8/20/12	ivee	3 A	
2012	8	8/20/12	ivee	3 B	
2012	8	8/20/12	ivee	3 B	
2012	8	8/20/12	ivee	3 B	
2012	8	8/20/12	ivee	3 B	
2012	8	8/20/12	ivee	3 C	
2012	8	8/20/12	ivee	3 C	
2012	8	8/20/12	ivee	3 C	
2012	8	8/20/12	ivee	3 D	
2012	8	8/20/12	ivee	4 B	
2012	8	8/20/12	ivee	5 A	
2012	8	8/20/12	ivee	5 A	
2012	8	8/20/12	ivee	5 A	64
2012	8	8/20/12	ivee	5 A	70
2012	8	8/20/12	ivee	5 B	58
2012	8	8/20/12	ivee	5 B	60
2012	8	8/20/12	ivee	5 B	60
2012	8	8/20/12	ivee	5 B	62
2012	8	8/20/12	ivee	5 B	86
2012	8	8/20/12	ivee	5 B	105
2012	8	8/20/12	ivee	6 B	57
2012	8	8/20/12	ivee	6 B	60
2012	8	8/20/12	ivee	6 B	300

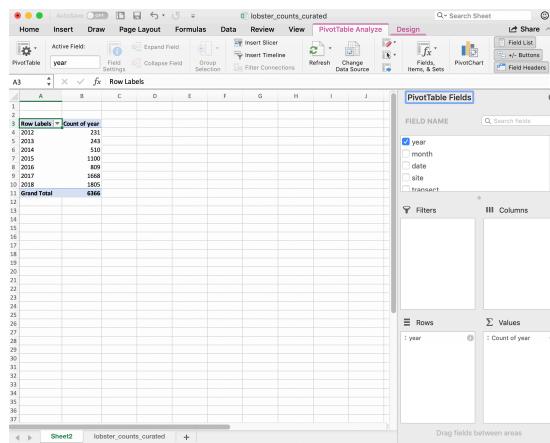
And it will do its best to find the data I would like to include in my Pivot Table (it can have difficulty with non-rectangular or “non-tidy” data), and suggest we make this in a new sheet:



And then we'll get a little wizard to help us create the Pivot Table.

### 6.6.1 pivot one variable

I want to summarize by year, so I drag “year” down into the “Rows” box, and to get the counts by year I actually drag the same variable, “year” into the “Values” box. And it will create a Pivot Table for me! But “sum” as the default summary statistic, so I can click the little “I” icon to change this to count.



A few things to note:

- The pivot table is separate entity from our data (it's on a different sheet); the original data has not been affected
- The pivot table only shows the variables we requested; we don't see other columns (like date, month, or site).
- Notice that in Excel we retain the overall totals for each site (in bold, on the same line with the site name). This is nice for communicating about data. But it can be problematic for further analyses, because it could be easy to take a total of this column and introduce errors.

So pivot tables are great because they summarize the data and keep the raw data raw — they even promote good practice because they by default ask you if you'd like to present the data in a new sheet rather than in the same sheet.

### 6.6.2 pivot two variables

We can also add site as a second variable by dragging it:

## 6.6. PIVOT TABLE DEMO

79

	A	B	C	D	E	F	G
3	Row Labels						
4	=Row#	=Count of year					
5	2012	33					
6	aqua	38					
7	carp	79					
8	date	26					
9	mook	83					
10	nopl	6					
11	=2013	243					
12	aqua	32					
13	carp	93					
14	date	40					
15	mook	15					
16	nopl	23					
17	=2014	510					
18	aqua	100					
19	carp	79					
20	date	132					
21	mook	36					
22	nopl	143					
23	=2015	1100					
24	aqua	83					
25	carp	90					
26	date	361					
27	mook	295					
28	nopl	270					
29	=2016	809					
30	aqua	45					
31	carp	231					
32	date	193					
33	mook	310					
34	nopl	127					
35	=2017	1668					
36	aqua	17					
37	carp	705					
38	date	606					
39	mook	178					
40	nopl	112					
41	=2018	1805					
42	aqua	54					
43	carp	343					
44	date	946					
45	mook	164					
46	nopl	298					
47	Grand Total	6366					

And then can reverse the order by dragging:

	A	B	C	D	E	F	G
3	Row Labels						
4	=Row#	=Count of year					
5	2012	32					
6	2013	32					
7	2014	100					
8	2015	83					
9	2016	48					
10	2017	97					
11	2018	54					
12	=carp	1619					
13	2012	78					
14	2013	93					
15	2014	79					
16	2015	90					
17	2016	231					
18	2017	705					
19	2018	343					
20	=ree	2304					
21	2012	26					
22	2013	40					
23	2014	132					
24	2015	941					
25	2016	193					
26	2017	606					
27	2018	946					
28	=mook	982					
29	2012	83					
30	2013	15					
31	2014	36					
32	2015	296					
33	2016	210					
34	2017	178					
35	2018	164					
36	=nopl	1039					
37	2012	6					
38	2013	63					
39	2014	163					
40	2015	270					
41	2016	127					
42	2017	112					
43	2018	298					
44	Grand Total	6366					

So in terms of our final interest of average size by site and year, we are on our way! I'm going to stop here because we want to be able to do this in R.

The power of R is in the automation, and in keeping that raw data truly raw.

Let's talk about how this looks like in R.

## 6.7 group\_by() %>% summarize()

In R, we can create the functionality of pivot tables by using 2 main `dplyr` verbs in combination: `group_by` and `summarize`.

Say it with me: “pivot tables are group\_by and then summarize”. And just like pivot tables, you have flexibility with how you are going to summarize. For example, we can calculate an average, or a total.

I think it’s incredibly powerful to visualize what we are talking about with our data when do do these kinds of operations. It looks like this (from RStudio’s cheatsheet; all cheatsheets available from <https://rstudio.com/resources/cheatsheets>):



When we were reporting by year or site, we were essentially modifying what we were grouping by (the different colors here in this figure).

Let’s do this in R.

### 6.7.1 group\_by one variable

Let’s try this on our `lobsters` data, just like we did in Excel. We will count the total number of lobster by year. In R vocabulary, we will group\_by year and then summarize by counting using `n()`, which is a function from `dplyr`. `n()` counts the number of times an observation shows up, and since this is uncounted data, this will count each row. We’ll also use the pipe operator `%>%`, which you can read as “and then”.

This to me reads: “take the lobsters data and then group\_by year and then summarize by count in a new column called ‘count’”

```
lobsters %>%
  group_by(year) %>%
  summarize(count_by_year = n())
```

```
## # A tibble: 7 x 2
##   year count_by_year
##   <dbl>        <int>
## 1 2012          231
## 2 2013          243
```

```
## # A tibble: 6 x 2
##   year     count
##   <dbl>     <dbl>
## 1 2014      510
## 2 2015     1100
## 3 2016      809
## 4 2017     1668
## 5 2018     1805
```

Notice how together, `group_by` and `summarize` minimize the amount of information we see. We also saw this with the pivot table. We lose the other columns that aren't involved here.

Question: What if you *don't* group\_by first? Let's try it and discuss what's going on.

```
lobsters %>%
  summarize(count = n())
```

```
## # A tibble: 1 x 1
##       count
##   <int>
## 1 6366
```

So if we don't `group_by` first, we will get a single summary statistic (sum in this case) for the whole dataset.

Another question: what if we *only* group\_by?

```
lobsters %>%
  group_by(year)
```

```
## # A tibble: 6,366 x 7
## # Groups:   year [7]
##   year month date    site transect replicate size_mm
##   <dbl> <dbl> <chr>  <chr>   <dbl> <chr>     <dbl>
## 1 2012     8 8/20/12 ivee      3 A         70
## 2 2012     8 8/20/12 ivee      3 B         60
## 3 2012     8 8/20/12 ivee      3 B         65
## 4 2012     8 8/20/12 ivee      3 B         70
## 5 2012     8 8/20/12 ivee      3 B         85
## 6 2012     8 8/20/12 ivee      3 C         60
## 7 2012     8 8/20/12 ivee      3 C         65
## 8 2012     8 8/20/12 ivee      3 C         67
## 9 2012     8 8/20/12 ivee      3 D         70
## 10 2012    8 8/20/12 ivee      4 B         85
## # ... with 6,356 more rows
```

### 6.7.2 RStudio Viewer

Let's now check the `lobsters` variable. We can do this by clicking on `lobsters` in the Environment pane in RStudio.

We see that we haven't changed any of our original data that was stored in this variable. (Just like how the pivot table didn't affect the raw data on the original sheet).

*Aside:* You'll also see that when you click on the variable name in the Environment pane, `View(lobsters)` shows up in your Console. `View()` (capital V) is the R function to view any variable in the viewer. So this is something that you can write in your RMarkdown script, although RMarkdown will not be able to knit this view feature into the formatted document. So, if you want include `View()` in your RMarkdown document you will need to either comment it out `#View()` or add `eval=FALSE` to the top of the code chunk so that the full line reads `{r, eval=FALSE}`.

### 6.7.3 group\_by multiple variables

Great. Now let's summarize by both year and site like we did in the pivot table. We are able to `group_by` more than one variable. Let's do this together:

```
lobsters %>%
  group_by(site, year) %>%
  summarize(count_by_siteyear = n())

## # A tibble: 35 x 3
## # Groups:   site [5]
##   site     year count_by_siteyear
##   <chr> <dbl>             <int>
## 1 aque    2012              38
## 2 aque    2013              32
## 3 aque    2014             100
## 4 aque    2015              83
## 5 aque    2016              48
## 6 aque    2017              67
## 7 aque    2018              54
## 8 carp    2012              78
## 9 carp    2013              93
## 10 carp   2014              79
## # ... with 25 more rows
```

text.

### 6.7.4 summarize multiple variables

We can summarize multiple variables at a time. So far we've done the count of lobster observations. Let's also do the mean and standard deviation. First let's use the `mean()` function to calculate the mean. We do this within the same `summarize()` function, but we can add a new line to make it easier to read. Notice how when you put your cursor within the parenthesis and hit return, the indentation will automatically align.

```
lobsters %>%
  group_by(site, year) %>%
  summarize(count_by_siteyear = n(),
            mean_size_mm = mean(size_mm))

## # A tibble: 35 x 4
## # Groups:   site [5]
##   site    year count_by_siteyear mean_size_mm
##   <chr> <dbl>           <int>        <dbl>
## 1 aque    2012            38          71
## 2 aque    2013            32         72.1
## 3 aque    2014           100         76.9
## 4 aque    2015            83         68.5
## 5 aque    2016            48         68.7
## 6 aque    2017            67         73.9
## 7 aque    2018            54         71.7
## 8 carp    2012            78         74.4
## 9 carp    2013            93         76.6
## 10 carp   2014            79          NA
## # ... with 25 more rows
```

*Aside* Command-I will properly indent selected lines.

Great! But this will actually calculate some of the means as NA because one or more values in that year are NA. So we can pass an argument that says to remove NAs first before calculating the average. Let's do that, and then also calculate the standard deviation with the `sd()` function:

```
lobsters %>%
  group_by(site, year) %>%
  summarize(count_by_siteyear = n(),
            mean_size_mm = mean(size_mm, na.rm=TRUE),
            sd_size_mm = sd(size_mm, na.rm=TRUE))

## # A tibble: 35 x 5
## # Groups:   site [5]
##   site    year count_by_siteyear mean_size_mm sd_size_mm
##   <chr> <dbl>           <int>        <dbl>       <dbl>
## 1 aque    2012            38          71         10.2
```

```

## 2 aque 2013      32      72.1    12.3
## 3 aque 2014     100      76.9    9.32
## 4 aque 2015      83      68.5    12.6
## 5 aque 2016      48      68.7    12.5
## 6 aque 2017      67      73.9    11.9
## 7 aque 2018      54      71.7    8.14
## 8 carp 2012      78      74.4    14.6
## 9 carp 2013      93      76.6    8.71
## 10 carp 2014     79      79.1    8.57
## # ... with 25 more rows

```

So we can make the equivalent of Excel's pivot table in R with `group_by` and then `summarize`. But a powerful thing about R is that maybe we want this information to be used in further analyses. We can make this easier for ourselves by saving this as a variable. So let's add a variable assignment to that first line:

```

siteyear_summary <- lobsters %>%
  group_by(site, year) %>%
  summarize(count_by_siteyear = n(),
            mean_size_mm = mean(size_mm, na.rm = TRUE),
            sd_size_mm = sd(size_mm, na.rm = TRUE))

```

`siteyear_summary`

```

## # A tibble: 35 x 5
## # Groups:   site [5]
##   site   year count_by_siteyear mean_size_mm sd_size_mm
##   <chr> <dbl>             <int>       <dbl>      <dbl>
## 1 aque   2012            38        71        10.2
## 2 aque   2013            32        72.1      12.3
## 3 aque   2014           100       76.9      9.32
## 4 aque   2015            83       68.5      12.6
## 5 aque   2016            48       68.7      12.5
## 6 aque   2017            67       73.9      11.9
## 7 aque   2018            54       71.7      8.14
## 8 carp   2012            78       74.4      14.6
## 9 carp   2013            93       76.6      8.71
## 10 carp  2014            79       79.1      8.57
## # ... with 25 more rows

```

### 6.7.5 Activity

1. Calculate the median `size_mm` (Hint: `?median`) and
2. create and `ggsave()` a plot.

Then, save, commit, and push your .Rmd, .html, and .png.

Solution (no peeking):

```
siteyear_summary <- lobsters %>%
  group_by(site, year) %>%
  summarize(count_by_siteyear = n(),
            mean_size_mm = mean(size_mm, na.rm = TRUE),
            sd_size_mm = sd(size_mm, na.rm = TRUE),
            median_size_mm = median(size_mm, na.rm = TRUE))

siteyear_summary

## a ggplot option:
ggplot(data = siteyear_summary, aes(x = year, y = median_size_mm, color = site)) +
  geom_line()
ggsave(here("figures", "lobsters-line.png"))

## another option:
ggplot(siteyear_summary, aes(x = year, y = median_size_mm)) +
  geom_col() +
  facet_wrap(~site)
ggsave(here("figures", "lobsters-col.png"))
```

Don't forget to knit, commit, and push!

Nice work everybody.

## 6.8 Oh no, our colleague sent the wrong data!

Oh no! After all our analyses and everything we've done, our colleague just emailed us at 4:30pm on Friday that he sent the wrong data and we need to redo all our analyses with a new .xlsx file: `lobsters2.xlsx`, not `lobsters.xlsx`. Aaaaah!

If we were doing this in Excel, this would be a bummer; we'd have to rebuild our pivot table and click through all of our logic again. And then export our figures and save them into our report.

But, since we did it in R, we are much safer. We can go back to the top of our RMarkdown file, and read in the updated dataset, and then re-knit. We will still need to check that everything outputs correctly, (and that column headers haven't been renamed), but our first pass will be to update the filename and re-knit:

```
## read in data
lobsters <- read_xlsx(here("data/lobsters2.xlsx"))
```

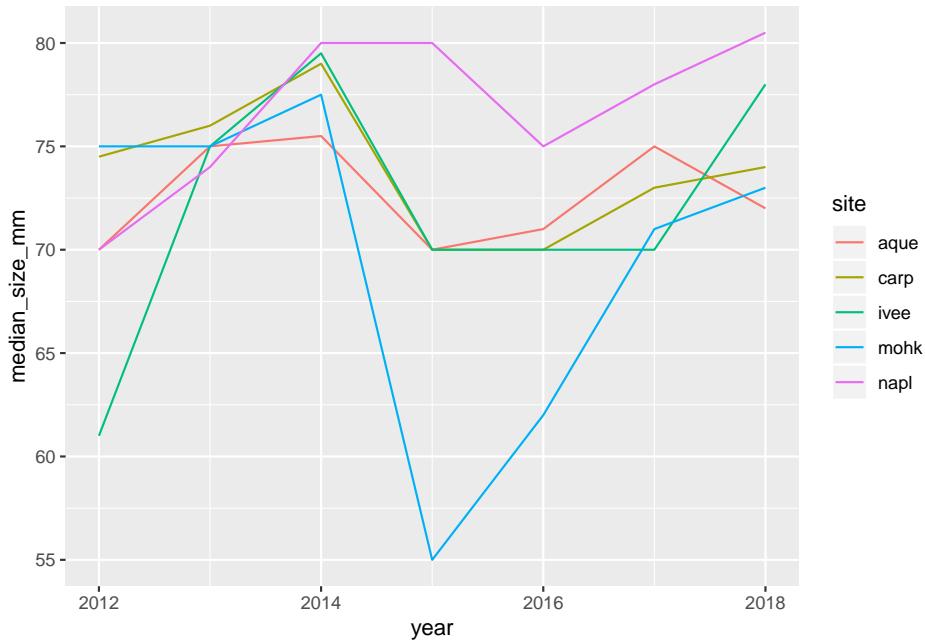
And now we can see that our plot updated as well:

```
siteyear_summary <- lobsters %>%
  group_by(site, year) %>%
  summarise(count_by_siteyear = n(),
            mean_size_mm = mean(size_mm, na.rm = TRUE),
            sd_size_mm = sd(size_mm, na.rm = TRUE),
            median_size_mm = median(size_mm, na.rm = TRUE), )

siteyear_summary
```

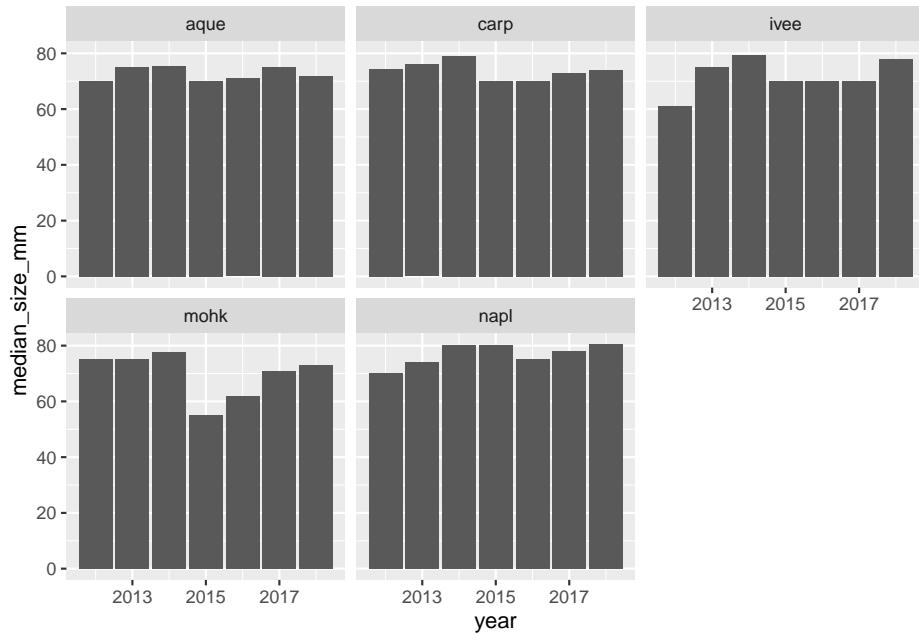
site	year	count_by_siteyear	mean_size_mm	sd_size_mm	median_size_mm
aque	2012	38	71	10.2	70
aque	2013	32	72.1	12.3	75
aque	2014	100	76.9	9.32	75.5
aque	2015	83	68.5	12.6	70
aque	2016	48	68.7	12.5	71
aque	2017	67	73.9	11.9	75
aque	2018	54	71.7	8.14	72
carp	2012	78	74.4	14.6	74.5
carp	2013	93	76.6	8.71	76
carp	2014	79	79.1	8.57	79

```
## # ... with 25 more rows
## a ggplot option:
ggplot(data = siteyear_summary, aes(x = year, y = median_size_mm, color = site)) +
  geom_line()
```



```
ggsave(here("figures", "lobsters-line.png"))
```

```
## Saving 6.5 x 4.5 in image
## another option:
ggplot(siteyear_summary, aes(x = year, y = median_size_mm)) +
  geom_col() +
  facet_wrap(~site)
```



```
ggsave(here("figures", "lobsters-col.png"))
```

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

### 6.8.1 Knit, push, & show differences on GitHub

So cool.

### 6.8.2 dplyr::count()

Now that we've spent time with group\_by %>% summarize, there is a shortcut if you only want to summarize by count. This is with a function called count(), and it will group\_by your selected variable, count, and then also ungroup. It looks like this:

```
lobsters %>%
  count(site, year)

## This is the same as:
lobsters %>%
  group_by(site, year) %>%
  summarise(n = n()) %>%
  ungroup()
```

Switching gears...

## 6.9 mutate()

### Make New Variables



There are a lot of times where you don't want to summarize your data, but you do want to operate beyond the original data. This is often done by adding a column. We do this with the `mutate()` function from `dplyr`. Let's try this with our original lobsters data. The sizes are in millimeters but let's say it was important for them to be in meters. We can add a column with this calculation:

```
# quick reminder what this looks like
head(lobsters)

## # A tibble: 6 x 7
##   year month date   site transect replicate size_mm
##   <dbl> <dbl> <chr> <chr>    <dbl> <chr>      <dbl>
## 1 2012     8 8/20/12 ivee       3 A          70
## 2 2012     8 8/20/12 ivee       3 B          60
## 3 2012     8 8/20/12 ivee       3 B          65
## 4 2012     8 8/20/12 ivee       3 B          70
## 5 2012     8 8/20/12 ivee       3 B          85
## 6 2012     8 8/20/12 ivee       3 C          60

lobsters %>%
  mutate(size_m = size_mm / 1000)

## # A tibble: 6,366 x 8
##   year month date   site transect replicate size_mm size_m
##   <dbl> <dbl> <chr> <chr>    <dbl> <chr>      <dbl>   <dbl>
## 1 2012     8 8/20/12 ivee       3 A          70  0.07
## 2 2012     8 8/20/12 ivee       3 B          60  0.06
## 3 2012     8 8/20/12 ivee       3 B          65  0.065
## 4 2012     8 8/20/12 ivee       3 B          70  0.07
## 5 2012     8 8/20/12 ivee       3 B          85  0.085
## 6 2012     8 8/20/12 ivee       3 C          60  0.06
## 7 2012     8 8/20/12 ivee       3 C          65  0.065
## 8 2012     8 8/20/12 ivee       3 C          67  0.067
## 9 2012     8 8/20/12 ivee       3 D          70  0.07
## 10 2012    8 8/20/12 ivee       4 B          85  0.085
```

```
## # ... with 6,356 more rows
```

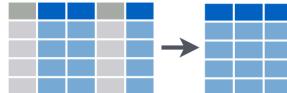
If we want to add a column that has the same value repeated, we can pass it just one value, either a number or a character string (in quotes). And let's save this as a variable called `lobsters_detailed`

```
lobsters_detailed <- lobsters %>%
  mutate(size_m = size_mm / 1000,
        millenia = 2000,
        observer = "Allison Horst")
```

## 6.10 `select()`

We will end with one final function, `select`. This is how to choose, retain, and move your data by columns:

### Subset Variables (Columns)



Let's say that we want to present this data finally with only columns for date, site, and size in meters. We would do this:

```
lobsters_detailed %>%
  select(date, site, size_m)

## # A tibble: 6,366 x 3
##   date     site   size_m
##   <chr>    <chr>   <dbl>
## 1 8/20/12  ivee    0.07
## 2 8/20/12  ivee    0.06
## 3 8/20/12  ivee    0.065
## 4 8/20/12  ivee    0.07
## 5 8/20/12  ivee    0.085
## 6 8/20/12  ivee    0.06
## 7 8/20/12  ivee    0.065
## 8 8/20/12  ivee    0.067
## 9 8/20/12  ivee    0.07
## 10 8/20/12  ivee   0.085
## # ... with 6,356 more rows
```

One last time, let's knit, save, commit, and push to GitHub.

## 6.11 Deep thoughts

Highly recommended read: Broman & Woo: Data organization in spreadsheets.  
Practical tips to make spreadsheets less error-prone, easier for computers to process, easier to share

Great opening line: “Spreadsheets, for all of their mundane rectangularness, have been the subject of angst and controversy for decades.”

## 6.12 Efficiency Tips

arrow keys with shift, option, command



# Chapter 7

## Tidying

### 7.1 Summary

In previous sessions, we learned to read in data, do some wrangling, and create a graph and table.

Here, we'll continue by *reshaping* data frames (converting from long-to-wide, or wide-to-long format), *separating* and *uniting* variable (column) contents, finding and replacing string patterns, and isolating numbers from number-character combinations.

### 7.2 Objectives

- Reshape data frames with `tidyverse::pivot_wider()` and `tidyverse::pivot_longer()`
- Convert column names with `janitor::clean_names()`
- Combine or separate information from columns with `tidyverse::unite()` and `tidyverse::separate()`
- Detect a string with `stringr::str_detect()`
- Replace a string with `stringr::str_replace()`
- Isolate numbers with `readr::parse_number()`
- Use our new skills as part of a bigger wrangling sequence

### 7.3 Resources

- Ch. 12 *Tidy Data*, in R for Data Science by Grolemund & Wickham - `tidyverse` documentation from tidyverse.org - `janitor` repo / information from Sam Firke

## 7.4 Set-up

### 7.4.1 Create a new R Markdown and attach packages

Within your day 2 R Project, create a new .Rmd. Attach the following packages (using `library(package_name)`):

- `tidyverse`
- `here`
- `janitor`
- `readxl`

Knit and save your new .Rmd within the project folder.

```
# Attach packages
library(tidyverse)
library(janitor)
library(here)
library(readxl)
```

### 7.4.2 Read in data with `read_excel()`

We've used both `read_csv()` and `read_excel()` to import data from spreadsheets into R.

Use `read_excel()` to read in the `inverts.xlsx` data as an object called `inverts`.

```
inverts <- read_excel(here("data", "inverts.xlsx"))
```

Be sure to explore the imported data a bit:

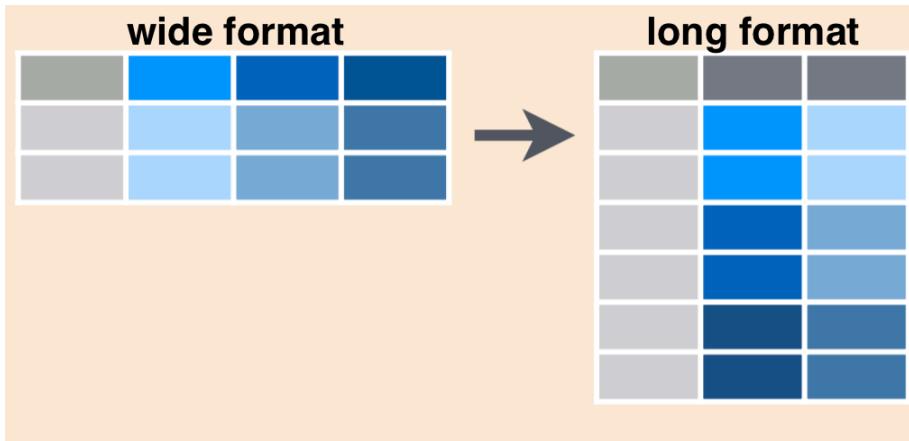
- `View()`
- `names()`
- `summary()`

## 7.5 Wide-to-longer format with `tidyr::pivot_longer()`

In *tidy format*, each variable is contained within a single column. If we look at `inverts_df`, we can see that the `year` variable is actually split over 3 columns, so we'd say this is currently in **wide format**.

There may be times when you want to have data in wide format, but often with code it is more efficient to convert to **long format** by gathering together observations for a variable that is currently split into multiple columns.

Schematically, converting from wide to long format looks like this:



We'll use `tidyverse::pivot_longer()` to gather data from all years in `inverts` (columns 2016, 2017, and 2018) into two columns:

- one called `year`, which contains the year (as a number)
- one called `sp_count` containing the number of each species observed.

The new data frame will be stored as `inverts_long`:

```
inverts_long <- pivot_longer(data = inverts,
                           cols = '2016':'2018',
                           names_to = "year",
                           values_to = "sp_count")
```

The outcome is the new long-format `inverts_long` data frame:

```
inverts_long
```

```
## # A tibble: 165 x 5
##   month site common_name      year  sp_count
##   <chr> <chr> <chr>        <chr>    <dbl>
## 1 7     abur  calif cone snail  2016     451
## 2 7     abur  calif cone snail  2017      28
## 3 7     abur  calif cone snail  2018     762
## 4 7     abur  calif spiny lobster 2016      17
## 5 7     abur  calif spiny lobster 2017      17
## 6 7     abur  calif spiny lobster 2018      16
## 7 7     abur  orange cup coral  2016      24
## 8 7     abur  orange cup coral  2017      24
## 9 7     abur  orange cup coral  2018      24
## 10 7    abur  purple urchin    2016      48
## # ... with 155 more rows
```

Hooray, long format!

One thing that isn't obvious at first (but would become obvious if you continued

working with this data) is that since those year numbers were initially column names (characters), when they are stacked into the *year* column, their class wasn't auto-updated to numeric.

Explore the class of *year* in *inverts\_long*:

```
class(inverts_long$year)
```

```
## [1] "character"
```

We'll use `dplyr::mutate()` in a different way here: to create a new column (that's how we've used `mutate()` previously) that has the same name of an existing column, in order to update and overwrite the existing column.

In this case, we'll `mutate()` to add a column called *year*, which contains an `as.numeric()` version of the existing *year* variable:

```
# Coerce "year" class to numeric:  
  
inverts_long <- inverts_long %>%  
  mutate(year = as.numeric(year))
```

Checking the class again, we see that *year* has been updated to a numeric variable:

```
class(inverts_long$year)
```

```
## [1] "numeric"
```

## 7.6 Long-to-wider format with `tidyverse::pivot_wider()`

In the previous example, we had information spread over multiple columns that we wanted to *gather*. Sometimes, we'll have data that we want to *spread* over multiple columns.

For example, imagine that starting from *inverts\_long* we want each species in the *common\_name* column to exist as its **own column**. In that case, we would be converting from a longer to a wider format, and will use `tidyverse::pivot_wider()`.

Specifically for our data, we'll use `pivot_wider()` to spread the *common\_name* across multiple columns as follows:

```
inverts_wide <- inverts_long %>%  
  pivot_wider(names_from = common_name,  
             values_from = sp_count)  
  
inverts_wide  
  
## # A tibble: 33 x 8
```

```

##   month site  year `california con~ `california spi~ `orange cup cor~
##   <chr> <chr> <dbl>           <dbl>           <dbl>           <dbl>
## 1 7    abur  2016        451          17          24
## 2 7    abur  2017        28           17          24
## 3 7    abur  2018       762          16          24
## 4 7    ahnd  2016        27           16          24
## 5 7    ahnd  2017        24           16          24
## 6 7    ahnd  2018        24           16          24
## 7 7    aque  2016      4971          48         1526
## 8 7    aque  2017      1752          48         1623
## 9 7    aque  2018      2616          48         1859
## 10 7   bull  2016     1735          24          36
## # ... with 23 more rows, and 2 more variables: `purple urchin` <dbl>, `rock
## #   scallop` <dbl>

```

We can see that now each *species* has its own column (wider format). But also notice that those column headers (since they have spaces) might not be in the most coder-friendly format...

## 7.7 Clean up column names with janitor::clean\_names()

The `janitor` package by Sam Firke is a brilliant collection of functions for some quick data cleaning. We recommend that you explore the different functions it contains. Like:

- `janitor::clean_names()`: update column headers to a case of your choosing
- `janitor::get_dups()`: see all rows that are duplicates within variables you choose
- `janitor::remove_empty()`: remove empty rows and/or columns
- `janitor::adorn_*`(): jazz up frequency tables of counts (we'll return to this for a table example in TODO: Session 8)

Here, we'll use `janitor::clean_names()` to convert all of our column headers to a more convenient case - the default is `lower_snake_case`, which means all spaces and symbols are replaced with an underscore (or a word describing the symbol), all characters are lowercase, and a few other nice adjustments.

For example, `janitor::clean_names()` would update these nightmare column names into much nicer forms:

- My...RECENT-income! becomes `my_recent_income`
- SAMPLE2.!test1 becomes `sample2_test1`
- ThisIsTheName becomes `this_is_the_name`
- 2015 becomes `x2015`

If we wanted to then use these columns (which we probably would, since we created them), we could clean the names to get them into more coder-friendly `lower_snake_case` with `janitor::clean_names()`:

```
inverts_wide <- inverts_wide %>%
  clean_names()

names(inverts_wide)

## [1] "month"                  "site"
## [3] "year"                   "california_cone_snail"
## [5] "california_spiny_lobster" "orange_cup_coral"
## [7] "purple_urchin"           "rock_scallop"
```

And there are other case options in `clean_names()`, like:

- “snake” produces `snake_case` (the default)
- “lower\_camel” or “small\_camel” produces `lowerCamel`
- “upper\_camel” or “big\_camel” produces `UpperCamel`
- “screaming\_snake” or “all\_caps” produces `ALL_CAPS`
- “lower\_upper” produces `lowerUPPER`
- “upper\_lower” produces `UPPERlower`

## 7.8 Combine or separate columns with `tidyr::unite()` and `tidyr::separate()`

Sometimes we’ll want to *separate* contents of a single column into multiple columns, or *combine* entries from different columns into a single column.

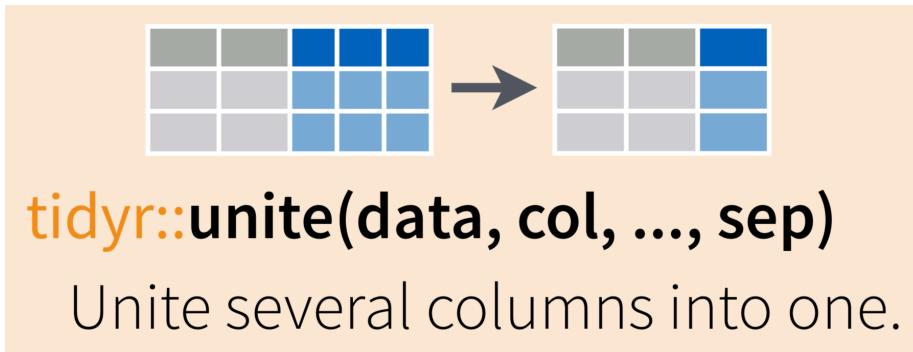
For example, the following data frame has *genus* and *species* in separate columns:

We may want to combine the genus and species into a single column, *scientific\_name*:

Or we may want to do the reverse (separate information from a single column into multiple columns). Here, we’ll learn `tidyr::unite()` and `tidyr::separate()` to help us do both.

### 7.8.1 `tidyr::unite()` to merge information from separate columns

Use `tidyr::unite()` to combine information from multiple columns into a single column (as for the scientific name example above)



To demonstrate uniting information from separate columns, we'll make a single column that has the combined information from *site* abbreviation and *year* in *inverts\_wide*.

We need to give `tidyr::unite()` several arguments:

- **data:** the data frame containing columns we want to combine (or pipe into the function from the data frame)
- **col:** the name of the new “united” column
- the **columns you are uniting**
- **sep:** the symbol, value or character to put between the united information from each column

```
inverts_unite <- inverts_wide %>%
  unite(col = "site_year", # What to name the new united column
        c(site, year), # The columns we'll unite (site, year)
        sep = "_") # How to separate the things we're uniting

## # A tibble: 6 x 7
##   month site_year california_cone~ california_spin~ orange_cup_coral
##   <chr> <chr>           <dbl>           <dbl>           <dbl>
## 1 7     abur_2016       451            17             24
## 2 7     abur_2017       28             17             24
## 3 7     abur_2018       762            16             24
## 4 7     ahnd_2016       27             16             24
## 5 7     ahnd_2017       24             16             24
## 6 7     ahnd_2018       24             16             24
## # ... with 2 more variables: purple_urchin <dbl>, rock_scallop <dbl>
```

Try updating the separator from “\_” to “hello!” to see what the outcome column contains.

`tidyr::unite()` can also combine information from *more* than two columns. For example, to combine the *site*, *common\_name* and *year* columns from *inverts\_long*, we could use:

```
# Uniting more than 2 columns:

inverts_triple_unite <- inverts_long %>%
  tidyverse::unite(col = "year_site_name",
                  c(year, site, common_name),
                  sep = "-")

head(inverts_triple_unite)

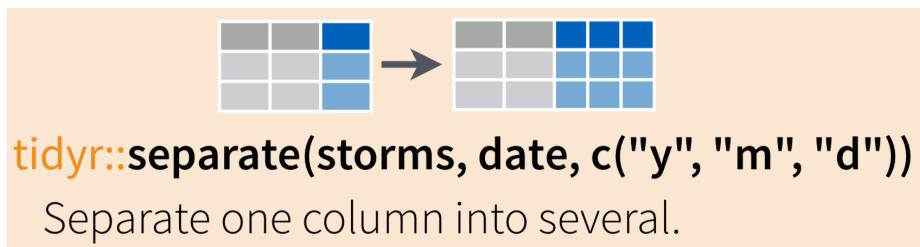
## # A tibble: 6 x 3
##   month year_site_name           sp_count
##   <chr>  <chr>                    <dbl>
## 1 7      2016-abur-california cone snail    451
## 2 7      2017-abur-california cone snail     28
## 3 7      2018-abur-california cone snail    762
## 4 7      2016-abur-california spiny lobster   17
## 5 7      2017-abur-california spiny lobster   17
## 6 7      2018-abur-california spiny lobster   16
```

### 7.8.2 `tidyverse::separate()` to separate information into multiple columns

While `tidyverse::unite()` allows us to combine information from multiple columns, it's more likely that you'll *start* with a single column that you want to split up into pieces.

For example, I might want to split up a column containing the *genus* and *species* (*Scorpaena guttata*) into two separate columns (*Scorpaena* | *guttata*), so that I can count how many *Scorpaena* organisms exist in my dataset at the genus level.

Use `tidyverse::separate()` to “separate a character column into multiple columns using a regular expression separator.”



Let's start again with `inverts_unite`, where we have combined the *site* and *year* into a single column called *site\_year*. If we want to **separate** those, we can use:

```
inverts_sep <- inverts_triple_unite %>%
  tidyverse::separate(year_site_name, into = c("my_year", "my_site_name"))

## Warning: Expected 2 pieces. Additional pieces discarded in 165 rows [1, 2, 3, 4,
## 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, ...].
```

What is that warning `Expected 2 pieces...` telling us? If we take a look at the resulting data frame `inverts_sep`, we see that it only keeps the first **two** pieces, and gets rid of the third (name). Which is a bit concerning, because we rarely want to just throw away information in a data frame.

```
head(inverts_sep)

## # A tibble: 6 x 4
##   month my_year my_site_name sp_count
##   <chr>  <chr>    <chr>        <dbl>
## 1 7      2016     abur         451
## 2 7      2017     abur          28
## 3 7      2018     abur         762
## 4 7      2016     abur          17
## 5 7      2017     abur          17
## 6 7      2018     abur          16
```

That's problematic. How can we make sure we're keeping as many different elements as exist in the united column?

We have a couple of options:

1. Create the *number* of columns that are needed to retain as many elements as exist (in this case, 3, but we only created two new columns in the example above)

```
inverts_sep3 <- inverts_triple_unite %>%
  tidyverse::separate(year_site_name, into = c("the_year", "the_site", "the_name"))

## Warning: Expected 3 pieces. Additional pieces discarded in 165 rows [1, 2, 3, 4,
## 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, ...].
```

Another warning. What is that about? Let's take a look at the resulting data frame and think about what's missing (what are the “pieces discarded”?):

```
head(inverts_sep3)

## # A tibble: 6 x 5
##   month the_year the_site the_name   sp_count
##   <chr>  <chr>    <chr>    <chr>        <dbl>
## 1 7      2016     abur     california     451
## 2 7      2017     abur     california     28
## 3 7      2018     abur     california    762
## 4 7      2016     abur     california     17
```

```
## 5 7      2017     abur      california      17
## 6 7      2018     abur      california      16
```

Aha! Only the *first word* of the common name was retained, and anything else was trashed. We want to keep everything after the second dash in the new *the\_name* column.

That's because the **default is extra = "warn"**, which means that if you have more pieces than columns you're separating into, it will populate the columns that have been allotted (in our case, just 3) then drop any additional information, giving you a warning that pieces have been dropped.

To keep the extra pieces that have been dropped, add the `extra = "merge"` argument within `tidy::separate()` to override:

```
inverts_sep_all <- inverts_triple_unite %>%
  separate(year_site_name,
           into = c("sample_year", "location", "sp_name"),
           extra = "merge")
```

No warning there about things being discarded. Explore *inverts\_sep\_all*:

```
## # A tibble: 165 x 5
##   month sample_year location sp_name          sp_count
##   <chr>    <chr>     <chr>   <chr>            <dbl>
## 1 7       2016       abur    california cone snail 451
## 2 7       2017       abur    california cone snail 28
## 3 7       2018       abur    california cone snail 762
## 4 7       2016       abur    california spiny lobster 17
## 5 7       2017       abur    california spiny lobster 17
## 6 7       2018       abur    california spiny lobster 16
## 7 7       2016       abur    orange cup coral    24
## 8 7       2017       abur    orange cup coral    24
## 9 7       2018       abur    orange cup coral    24
## 10 7      2016       abur    purple urchin     48
## # ... with 155 more rows
```

We see that the resulting data frame has split *year\_site\_name* into three separate columns, *sample\_year*, *location*, and *sp\_name*, but now everything after the second break ("–") remains together in *sp\_name* instead of dropping pieces following the third word.

## 7.9 Get just the numbers with `readr::parse_number()`

Sometimes, data are entered in a way that combines numeric information with other character bits. Like:

- 152.4 grams

## 7.10. DETECTING A STRING PATTERN WITH `STRINGR::STR_DETECT()` 103

- \$ 800.23
- 54.6 %
- mass 458.9 lbs

But usually, we want to work with the numeric piece on its own.

**Note:** keep that in mind for data entry - it's good to include units, symbols, and descriptive characters, but include them in a separate column, and/or in metadata!

We can isolate *just* the numeric part of an entry using `readr::parse_number()`.

As an example, let's return to the `inverts_unite` object we created previously. Let's say that we want to create a new column that contains only the numeric part from `site_year`.

We could use `parse_number()` as follows:

```
parse_yr <- inverts_unite %>%
  mutate(
    year_num = parse_number(site_year)
  )
```

Notice that `parse_yr` now has an additional column, `year`, which only contains the numeric piece from the `site_year` column. Cool!

## 7.10 Detecting a string pattern with `stringr::str_detect()`

Sometimes we'll want to find observations that contain a specific string pattern within a variable of interest.

For example, consider the fantasy data below:

There might be a time when we would want to use observations that:

- Contain the string “fish”, in isolation or within a larger string (like “rock-fish”)
- Contain the string “blue”

In those cases, it would be useful to `detect` a string pattern. Here, we'll use `stringr::str_detect()` to find and keep observations that contain our specified string pattern.

Refamiliarize yourself with the `inverts` object:

```
View(inverts)
```

First, we want to detect observations where the `site` variable contains string pattern “sc”. Note that there are two sites that contain the pattern: `sedi` and `sctw`.

Using `str_detect()` to find and keep observations where the `site` variable contains pattern “sc”:

```
sc_detect <- inverts %>%
  filter(
    str_detect(site, pattern = "sc")
  )

# Use `unique` to return the unique entries in the site column:
unique(sc_detect$site)

## [1] "scdi" "sctw"
```

We see that only sites `scdi` and `sctw` remain in `sc_detect`.

### 7.10.1 Activity

Create a new object called `in_detect`, starting from `inverts`, that only contains observations if the `common_name` variable contains the string pattern “in”.

What species remain?

Solution:

```
in_detect <- inverts %>%
  filter(
    str_detect(common_name, pattern = "in")
  )

unique(in_detect$common_name)

## [1] "california spiny lobster" "purple urchin"
# Only 'california spiny lobster' and 'purple urchin' remain!
```

We can similarly choose to *exclude* observations that contain a set string pattern by adding the `negate = TRUE` argument within `str_detect()`.

For example, to exclude any observations where the common name contains “california,” we’d use:

```
no_ca <- inverts %>%
  filter(
    str_detect(common_name, pattern = "california", negate = TRUE)
  )
```

Look at `no_ca` to confirm that “california spiny lobster” and “california cone snail” have been excluded.

## 7.11 Replacing a string with `stringr::str_replace()`

Was data entered in a way that's difficult to code with, or is just plain annoying? Did someone wrongly enter “fish” as “fsh” throughout the spreadsheet, and you want to update it everywhere?

Use `stringr::str_replace()` to automatically replace a string pattern.

**Warning:** The pattern will be replaced everywhere - so if you ask to replace “fsh” with “fish”, then “offshore” would be updated to “offishore”. Be careful to ensure that when you think you’re making one replacement, you’re not also replacing something else unexpectedly.

Starting with `inverts`, let’s any place we find “california” we want to replace it with the abbreviation “CA”:

```
ca_abbr <- inverts %>%
  mutate(
    common_name =
      str_replace(common_name,
                  pattern = "california",
                  replacement = "CA")
  )
```

Now, check to confirm that “california” has been replaced with “CA”.

### 7.11.1 END tidying session!



# Chapter 8

## Dplyr and vlookups

### 8.1 Summary

In previous sessions, we've learned to do some basic wrangling and find summary information with functions in the `dplyr` package, which exists within the `tidyverse`. We've used:

- `count()`: get counts of observations for groupings we specify
- `mutate()`: add a new column, while keeping the existing ones
- `group_by()`: let R know that `groups` exist within the dataset, by variable(s)
- `summarize()`: calculate a value (that you specify) for each group, then report each group's value in a table

In this session, we'll expand our data wrangling toolkit using:

- `filter()` to conditionally subset our data by `rows`, and
- `*_join()` functions to merge data frames together
- Make a nicely formatted table with `kable()` and `kableExtra`

The combination of `filter()` and `*_join()` - to return rows satisfying a condition we specify, and merging data frames by like variables - is analogous to the useful VLOOKUP function in Excel.

#### 8.1.1 Objectives

- Use `filter()` to subset data frames, returning `rows` that satisfy variable conditions
- Use `full_join()`, `left_join()`, and `inner_join()` to merge data frames, with different endpoints in mind
- Use `filter()` and `*_join()` as part of a wrangling sequence

### 8.1.2 Resources

- `filter()` documentation from tidyverse.org
- `join()` documentation from tidyverse.org
- Chapters 5 and 13 in *R for Data Science* by Garrett Grolemund and Hadley Wickham
- “Create awesome HTML tables with `knitr::kable()` and `kableExtra`” by Hao Zhu

## 8.2 Set-up: Create a new .Rmd, attach packages & get data

Create a new R Markdown document in your r-workshop project and knit to save as **filter\_join.Rmd**. Remove all the example code (everything below the set-up code chunk).

In this session, we’ll use four packages:

- `tidyverse`
- `readxl`
- `here`
- `kableExtra`

Attach the packages in the setup code chunk in your .Rmd:

```
# Attach packages:
library(tidyverse)
library(readxl)
library(here)
library(kableExtra)
```

Then create a new code chunk to read in three files from your ‘data’ subfolder:

- `inverts.xlsx`
- `fish.csv`
- `kelp.xlsx`

```
# Read in data:
inverts <- read_excel(here("data", "inverts.xlsx"))
fish <- read_csv(here("data", "fish.csv"))
kelp <- read_excel(here("data", "kelp.xlsx"))
```

We should always explore the data we’ve read in using functions. Use `View()`, `names()`, `summary()`, `head()` and `tail()` to check it out.

Now, let’s use `filter()` to decide which observations (rows) we’ll keep or exclude in new subsets, similar to using Excel’s VLOOKUP function.

## 8.3 dplyr::filter() to conditionally subset by rows

Use `filter()` to let R know which **rows** you want to keep or exclude, based whether or not their contents match conditions that you set for one or more variables.

### Subset Observations (Rows)



Some examples in words that might inspire you to use `filter()`:

- “I only want to keep rows where the temperature is greater than 90°F.”
- “I want to keep all observations **except** those where the tree type is listed as **unknown**.”
- “I want to make a new subset with only data for mountain lions (the species variable) in California (the state variable).”

When we use `filter()`, we need to let R know a couple of things:

- What data frame we’re filtering from
- What condition(s) we want observations to **match** and/or **not match** in order to keep them in the new subset

Here, we’ll learn some common ways to use `filter()`.

#### 8.3.0.1 Filter rows by matching a single character string

Let’s say we want to keep all observations from the **fish** data frame where the common name is “garibaldi.” Here, we need to tell R to only *keep rows* from the **fish** data frame when the common name (`common_name` variable) exactly matches **garibaldi**. Use `==` to ask R to look for matching strings:

```
fish_garibaldi <- filter(fish, common_name == "garibaldi")
```

Check out the **fish\_garibaldi** object to ensure that only *garibaldi* observations remain.

You could also do this using the pipe operator `%>%`:

```
fish_garibaldi <- fish %>%
  filter(common_name == "garibaldi")
```

### 8.3.1 Activity

**Task:** Create a subset from the **fish** data frame, stored as object **fish\_abur**, that only contains observations from Arroyo Burro (site ‘abur’).

**Solution:**

```
fish_abur <- fish %>%
  filter(site == "abur")
```

Explore the subset you just created to ensure that only Arroyo Burro observations are returned.

### 8.3.2 Filter rows based on numeric conditions

Use expected operators ( $>$ ,  $<$ ,  $\geq$ ,  $\leq$ ,  $==$ ) to set conditions for a numeric variable when filtering. For this example, we only want to retain observations when the **total\_count** column value is  $\geq 50$ :

```
fish_over50 <- filter(fish, total_count >= 50)
```

Or, using the pipe:

```
fish_over50 <- fish %>%
  filter(total_count >= 50)
```

### 8.3.3 Filter to return rows that match *this OR that OR that*

What if we want to return a subset of the **fish** df that contains *garibaldi*, *blacksmith* OR *black surfperch*?

There are several ways to write an “OR” statement for filtering, which will keep any observations that match Condition A *or* Condition B *or* Condition C. In this example, we will create a subset from **fish** that only contains rows where the **common\_name** is *garibaldi* or *blacksmith* or *black surfperch*.

Use **%in%** to ask R to look for *any matches* within a combined vector of strings:

```
fish_3sp <- fish %>%
  filter(common_name %in% c("garibaldi", "blacksmith", "black surfperch"))
```

Alternatively, you can indicate **OR** using the vertical line operator **|** to do the same thing (but you can see that it’s more repetitive when looking for matches within the same variable):

```
fish_3sp <- fish %>%
  filter(common_name == "garibaldi" | common_name == "blacksmith" | common_name == "bl
```

### 8.3.4 Filter to return rows that match conditions for multiple variables

In the previous examples, we set filter conditions based on a single variable (e.g. `common_name`). What if we want to return observations that satisfy conditions for multiple variables?

For example: We want to create a subset that only returns rows from `inverts` where the `site` is “abur” or “mohk” *and* the `common_name` is “purple urchin.” In `filter()`, add a comma (or ampersand ‘&’) between arguments for multiple “and” conditions:

```
urchin_abur_mohk <- inverts %>%
  filter(site %in% c("abur", "mohk"), common_name == "purple urchin")

head(urchin_abur_mohk)

## # A tibble: 2 x 6
##   month site  common_name    `2016` `2017` `2018`
##   <chr> <chr> <chr>       <dbl>   <dbl>   <dbl>
## 1 7     abur  purple urchin     48     48     48
## 2 7     mohk  purple urchin    620    505    323
```

Like most things in R, there are other ways to do the same thing. For example, you could do the same thing using `&` (instead of a comma) between “and” conditions:

```
# Use the ampersand (&) to add another condition "and this must be true":

urchin_abur_mohk <- inverts %>%
  filter(site %in% c("abur", "mohk") & common_name == "purple urchin")
```

Or you could just do two filter steps in sequence:

```
# Written as sequential filter steps:

urchin_abur_mohk <- inverts %>%
  filter(site %in% c("abur", "mohk")) %>%
  filter(common_name == "purple urchin")
```

### 8.3.5 Activity: combined filter conditions

**Task:** Create a subset from the `fish` data frame, called `low_gb_wr` that only contains:

- Observations of *garibaldi* and *rock wrasse*
- Where the `total_count` is *less than or equal to 10*

**Solution:**

```
low_gb_wr <- fish %>%
  filter(common_name %in% c("garibaldi", "rock wrasse"),
         total_count <= 10)
```

Sync your local project to your repo on GitHub.

### 8.3.6 Example: combining `filter()` with other functions using the pipe operator (`%>%`)

We can also use `filter()` in combination with the functions we previously learned for wrangling. If we have multiple sequential steps to perform, we can string them together using the *pipe operator* (`%>%`).

Here, we'll start with the `inverts` data frame and create a subset that:

- Converts to long(er) format with `pivot_longer()`
- Only keeps observations for rock scallops
- Calculates the total count of rock scallops by site only

```
# Counts of scallops by site (all years included):
scallop_count_by_site <- inverts %>%
  pivot_longer(cols = '2016':'2018',
               names_to = "year",
               values_to = "sp_count") %>%
  filter(common_name == "rock scallop") %>%
  group_by(site) %>%
  summarize(tot_count = sum(sp_count, na.rm = TRUE))

scallop_count_by_site

## # A tibble: 11 x 2
##   site    tot_count
##   <chr>     <dbl>
## 1 abur      48
## 2 ahnd      48
## 3 aque     152
## 4 bull      48
## 5 carp    2519
## 6 golb      48
## 7 ivee     169
## 8 mohk      346
## 9 napl    6416
## 10 scdi     2390
## 11 sctw     1259
```

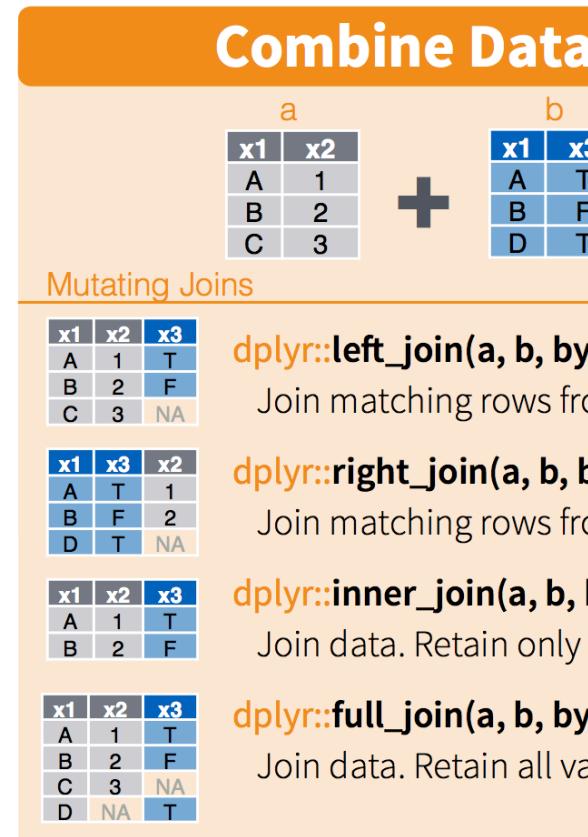
## 8.4 Join data frames with dplyr::\*\_join()

Excel's VLOOKUP can also be used to merge data from separate tables or worksheets. Here, we'll use the `*_join()` functions to merge separate data frames in R.

There are a number of ways to merge data frames in R. We'll use `full_join()`, `left_join()`, and `inner_join()` in this session.

From R Documentation (`?join`):

- `full_join()`: “returns all rows and all columns from both x and y. Where there are not matching values, returns NA for the one missing.” Basically, nothing gets thrown out, even if a match doesn't exist - making `full_join()` the safest option for merging data frames. When in doubt, `full_join()`.
- `left_join()`: “return all rows from x, and all columns from x and y. Rows in x with no match in y will have NA values in the new columns. If there are multiple matches between x and y, all combinations of the matches are returned.”
- `inner_join()`: “returns all rows from x where there are matching values in y, and all columns from x and y. If there are multiple matches between x and y, all combination of the matches are returned.” This will drop observations that don't have a match between the merged data frames, which makes it a riskier merging option if you're not sure what you're trying to do.



Schematic (from RStudio data wrangling cheat sheet):

To clarify what the different joins are doing, let's first make a subset of the *fish* data frame that only contains observations from 2016 and 2017.

```
fish_2016_2017 <- fish %>%
  filter(year == 2016 | year == 2017)
```

Take a look to ensure that only those years are included with `View(fish_2016_2017)`. Now, let's merge it with our Arroyo Burro kelp fronds data (`kelp_abur`) in different ways.

### 8.4.1 `full_join()` to merge data frames, keeping everything

When we join data frames in R, we need to tell R a couple of things (and it does the hard joining work for us):

- Which data frames we want to merge together
- Which variables to merge by

**Note:** If there are **exactly matching** column names in the data frames you're merging, the `*_join()` functions will assume that you want to join by those columns. If there are *no* matching column names, you can specify which columns to join by manually. We'll do both here.

```
# Join the fish_2016_2017 and kelp_abur
abur_kelp_join <- fish_2016_2017 %>%
  full_join(kelp_abur, by = c("year", "site")) # Uh oh. An error message.
```

When we try to do that join, we get an error message: `Error: Can't join on 'year' x 'year' because of incompatible types (character / numeric)`

Let's google this. That means copying this from the console and pasting it into Google.

What's going on here? First, there's something fishy (ha) going on with the class of the `year` variable in `kelp_abur`. Use the `class()` function to see how R understands that variable (remember, we use `$` to return a specific column from a data frame).

```
class(kelp_abur$year)
```

```
## [1] "character"
```

So the variable is currently stored as a character. Why?

If we go back to the `kelp.xlsx` file, we'll see that the numbers in both the year and month column have been stored as *text*. There are several hints Excel gives us:

- Cells are left aligned, when values stored as numbers are right aligned
- The green triangles in the corner indicate some formatting
- The warning sign shows up when you click on one of the values with text formatting, and lets you know that the cell has been stored as text. We are given the option to reformat as numeric in Excel, but we'll do it here in R so we have a reproducible record of the change to the variable class.

There are a number of ways to do this in R. We'll use `mutate()` to overwrite the existing `year` column while coercing it to class *numeric* using the `as.numeric()` function.

```
# Coerce the class of 'year' to numeric
kelp_abur <- kelp_abur %>%
  mutate(year = as.numeric(year))
```

Now if we check the class of the `year` variable in `kelp_counts_abur`, we'll see that it has been coerced to 'numeric':

```
class(kelp_abur$year)
```

```
## [1] "numeric"
```

**Question: Isn't it bad practice to overwrite variables, instead of just making a new one?** Great question, and usually the answer is yes. Here, we feel fine with “overwriting” the year column because we’re not changing anything about what’s contained within the column, we’re only changing how R understands it. Always use caution if overwriting variables, and if in doubt, add one instead!

OK, so now the class of *year* in the data frames we’re joining is the same. Let’s try that `full_join()` again:

```
abur_kelp_join <- fish_2016_2017 %>%
  full_join(kelp_abur, by = c("year", "site"))
```

Let’s look at the merged data frame with `View(abur_kelp_join)`. A few things to notice about how `full_join()` has worked:

1. All columns that existed in **both data frames** still exist.
2. All observations are retained, even if they don’t have a match. In this case, notice that for other sites (not ‘abur’) the observation for fish still exists, even though there was no corresponding kelp data to merge with it. The kelp frond data from 2018 is also returned, even though the fish counts dataset did not have ‘year == 2018’ in it.
3. The kelp frond data is joined to *all observations* where the joining variables (*year*, *site*) are a match, which is why it is repeated 5 times for each year (once for each fish species).

Because all data (observations & columns) are retained, `full_join()` is the safest option if you’re unclear about how to merge data frames.

#### 8.4.2 `left_join()` to merge data frames, keeping everything in the ‘x’ data frame and only matches from the ‘y’ data frame

Now, we want to keep all observations in *fish\_2016\_2017*, and merge them with *kelp\_abur* while only keeping observations from *kelp\_abur* that match an observation within *fish\_2016\_2017*. So when we use `left_join()`, any information on kelp frond counts from 2018 should be dropped, because those wouldn’t have a match in the left data frame.

```
fish_kelp_2016_2017 <- fish_2016_2017 %>%
  left_join(kelp_abur)
```

```
## Joining, by = c("year", "site")
```

Notice when you look at *fish\_kelp\_2016\_2017*, the 2018 data that **does** exist in *kelp\_abur* does **not** get joined to the *fish\_2016\_2017* data frame, because `left_join(df_a, df_b)` will only keep observations from *df\_b* if they have a match in *df\_a*!

### 8.4.3 `inner_join()` to merge data frames, only keeping observations with a match in both

Use `inner_join()` if you **only** want to retain observations that have matches across **both data frames**. Caution: this is built to exclude any observations that don't match across data frames by joined variables - double check to make sure this is actually what you want to do!

For example, if we use `inner_join()` to merge `fish` and `kelp_abur`, then we are asking R to **only return observations where the joining variables (`year` and `site`) have matches in both data frames**. Let's see what the outcome is:

```
kelp_fish_injoin <- fish %>%
  inner_join(kelp_abur)

## Joining, by = c("year", "site")
# kelp_fish_injoin
```

Here, we see that only observations (rows) where there is a match for `year` and `site` in both data frames are returned.

### 8.4.4 Using `filter()` and `join()` in a sequence

Now let's combine what we've learned about piping, filtering and joining!

Let's complete the following as part of a single sequence (remember, check to see what you've produced after each step) to create a new data frame called `my_fish_join`:

- Start with **fish** data frame
- Filter **fish** to only including observations for 2017 at Arroyo Burro
- Join the **kelp\_abur** data frame to the resulting subset using `left_join()`
- Add a new column that contains the ‘fish per kelp fronds’ density (total\_count / total\_fronds)

That sequence might look like this:

```
my_fish_join <- fish %>%
  filter(year == 2017, site == "abur") %>%
  left_join(kelp_abur, by = c("year", "site")) %>%
  mutate(fish_per_frond = total_count / total_fronds)
```

Explore the resulting **my\_fish\_join** data frame.

## 8.5 A nice HTML table with `kable()` and `kableExtra`

With any data frame, you can a nicer looking table in your knitted HTML using `knitr::kable()` and functions in the `kableExtra` package.

Start by using `kable()` with `my_fish_join`, and see what the default HTML table looks like in your knitted document:

```
kable(kelp_abur)
```

Simple, but quick to get a clear & useful table! Now let's spruce it up a bit with `kableExtra::kable_styling()` to modify HTML table styles:

```
kelp_abur %>%
  kable() %>%
  kable_styling(bootstrap_options = "striped", full_width = FALSE)
```

...with many other options for customizing HTML tables! Make sure to check out “Create awesome HTML tables with `knitr::kable()` and `kableExtra`” by Hao Zhu for more examples and options.

We'll make several more tables later on in the workshop.

### 8.5.1 Sync your project with your repo on GitHub

## 8.6 Fun / kind of scary facts

### How is this similar to VLOOKUP in Excel? How does it differ?

From Microsoft Office Support, “use VLOOKUP when you need to find things in a table or a range by row.”

So, both `filter()` and `VLOOKUP` look through your data frame (or spreadsheet, in Excel) to look for observations that match your conditions. But they also differ in important ways:

- (1) By default `VLOOKUP` looks for and returns an observation for *approximate* matches (and you have to change the final argument to `FALSE` to look for an exact match). In contrast, by default `filter()` will look for exact conditional matches.
- (2) `VLOOKUP` will look for and return information from the *first observation* that matches (or approximately matches) a condition. `filter()` will return all observations (rows) that exactly match a condition.

## 8.7 Efficiency Tips

- Comment out multiline code with Command + Shift + C
- Knit with Command + Shift + K

## 8.8 End `filter()` + `join()` section



# Chapter 9

## Collaborating

### 9.1 Summary

#### 9.1.1 Objectives

- create a new repo and give permission to a collaborator
- open as a new RStudio project, collaborate with a partner
- explore github.com blame, history, issues
- How to effectively ask for help
  - R communities and getting help (rOpenSci, RLadies)
  - reprex & RStudio Community “do it with iris”
  - how to use Twitter for #rststats Activity: reprex to GitHub “dplyr reprex”

#### 9.1.2 Resources

- Allison Horst ESM 206 Lecture 2

### 9.2 Lesson

### 9.3 Troubleshooting: help you help yourself

Borrow/show Allison’s slides - Allison Horst ESM 206 Lecture 2

### 9.3.1 Error messages are your friends

## 9.4 Efficiency Tips

## 9.5 Additional thoughts

---

always\_allow\_html: true

---

# Chapter 10

## Synthesis

### 10.1 Summary

In this session, we'll pull together the skills that we've learned so far. We'll create a new GitHub repo and R project, wrangle and visualize data from spreadsheets in R Markdown, communicate between RStudio (locally) and GitHub (remotely) to keep our updates safe, then TODO: share our outputs in a nicely formatted GitHub page. And we'll learn a few new things along the way!

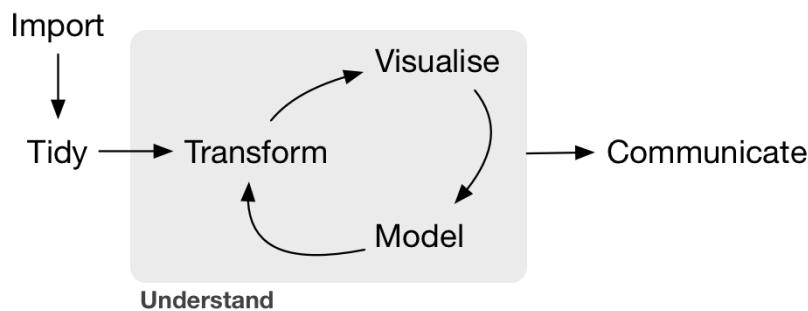


Figure 10.1: Grolemund & Wickham R4DS Illustration

### 10.2 Objectives

- Create a new repo on GitHub
- Clone the repo to create a new version-controlled R project
- Add data to \*\*data\*\* folder in working directory

- Create a new R Markdown document
- Attach necessary packages (`tidyverse`, `here`, `janitor`, `kableExtra`)
- Use `here::here()` for simpler (and safer) file paths
- Data tidying and wrangling (`dplyr`, `tidyr`, etc.)
- HTML tables with `kable()` and `kableExtra`
- Data visualization (`ggplot2`)
- TODO: Publish with a useful ReadMe to share?

### 10.3 Resources

- Project oriented workflows by Jenny Bryan

### 10.4 Set-up

- In GitHub, create a new repository called `us-fisheries-landings`
- Clone the repo to create a version controlled project (remember, copy & paste the URL from the GitHub Clone / Download)
- In the local project working directory, create a subfolder called ‘data’
- Copy and paste the `noaa_fisheries.csv` file into the ‘data’ folder
- Create a new R Markdown document within your `us-fisheries-landings` project
- Knit your `.Rmd` to `html`, saving as `us_landings.Rmd`

**Data:** File name: `noaa_fisheries.csv` Description: NOAA Commercial Fisheries Landing data (1950 - 2017) Accessed from: <https://www.st.nmfs.noaa.gov/commercial-fisheries/commercial-landings/> Source: Fisheries Statistics Division of the NOAA Fisheries

*Note on the data:* “aggregate” here means “These names represent aggregations of more than one species. They are not inclusive, but rather represent landings where we do not have species-specific data. Selecting “Sharks”, for example, will not return all sharks but only those where we do not have more specific information.”

### 10.5 Attach packages, read in and explore the data

Attach (load) packages with `library()`:

```
library(tidyverse)
library(here)
```

```
library(janitor)
library(kableExtra)
```

Read in the noaa\_fisheries.csv data as object `us_landings`:

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

```
## Parsed with column specification:
## cols(
##   Year = col_double(),
##   State = col_character(),
##   `AFS Name` = col_character(),
##   `Landings (pounds)` = col_double(),
##   `Dollars (USD)` = col_character()
## )
```

Go exploring a bit:

```
summary(us_landings)
View(us_landings)
names(us_landings)
head(us_landings)
tail(us_landings)
```

## 10.6 Some data cleaning to get salmon landings

- Clean up column names (`clean_names()`)
- Convert everything to lower case with `mutate() + (str_to_lower())`
- Remove dollar signs in value column (`mutate() + parse_number()`)
- Remove the “aggregate” grouping indicator in species (`mutate() + a new function! str_remove() %>%`)
- Keep only observations that include “salmon” (`filter()`)
- Separate the grouped name (“salmon”) from any additional refined information on species (`separate()`)

We'll break this up into two pieces:

1. A tidier version of the entire data frame
2. A subset that only contains salmon information

First: tidying the entire data frame

```
landings_tidy <- us_landings %>%
  clean_names() %>% # Make column headers snake_case
  mutate(
    state = str_to_lower(state),
    afs_name = tolower(afs_name)
```

```
) %>% # Converts character columns to lowercase
mutate(dollars_num = parse_number(dollars_usd))
```

Now, getting just the salmon:

```
salmon_landings <- landings_tidy %>%
  mutate(afs_clean = str_remove(afs_name, pattern = "aggregate")) %>%
  filter(str_detect(afs_clean, pattern = "salmon")) %>% # Detect matches for crab
  separate(afs_clean, into = c("group", "species"), sep = ",")
```

Explore `salmon_landings`.

## 10.7 Find grouped summary data

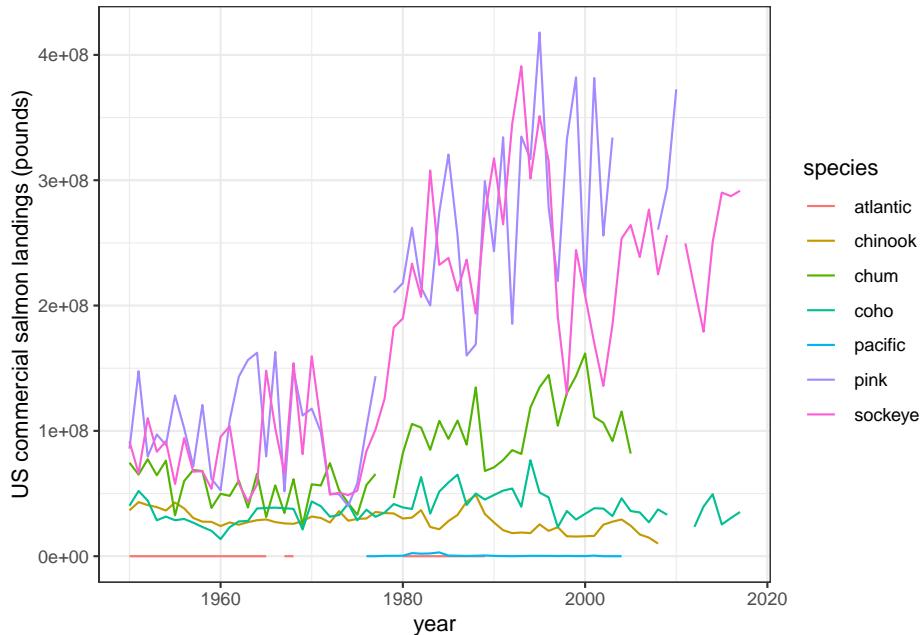
Find the annual total US landings and dollar value (summing across all states) for each type of salmon using `group_by()` + `summarize()`:

```
salmon_summary <- salmon_landings %>%
  group_by(year, species) %>%
  summarize(
    tot_landings = sum(landings_pounds),
    tot_value = sum(dollars_num)
  )
```

## 10.8 Make a graph of US commercial fisheries value by species over time with ggplot2

```
salmon_landings_graph <- ggplot(salmon_summary, aes(x = year, y = tot_landings, group =
  geom_line(aes(color = species)) +
  theme_bw() +
  labs(x = "year", y = "US commercial salmon landings (pounds)")

salmon_landings_graph
```



## 10.9 Export your salmon value graph with `ggsave`

```
ggsave(plot = salmon_landings_graph, here("figures", "us_salmon_landings.png"),
       height = 5, width = 8)
```

## 10.10 2015 commercial fisheries value by state

Now, let's create a finalized table of the top 5 states (by total commercial fisheries value) for 2015 .

Remember that we already created a tidied data frame, `landings_tidy`.

*Critical thinking for data wrangling workflow: Why does it make sense (especially now) that we broke our previous wrangling into two steps before getting our salmon subset?*

```
state_value <- landings_tidy %>%
  filter(year %in% c(2015)) %>%
  group_by(state) %>%
  summarize(
    state_value = sum(dollars_num, na.rm = TRUE),
```

```
state_landings = sum(landings_pounds, na.rm = TRUE)
) %>%
arrange(-state_value) %>%
head(5)
```

Look at our summary data frame, `state_value` - which contains all the information we want, but doesn't look very finalized.

What are some ways we'd want to finalize it for a report?

- Change units (to millions of dollars and pounds)
- Update column names
- Capitalize the state names

Let's do it!

## 10.11 Making a nice HTML table

First, we'll create it as a finalized data frame:

```
state_table <- state_value %>%
  mutate(`Fisheries value ($ millions)` = round(state_value / 1e6, 2),
        `Landings (million pounds)` = round(state_landings / 1e6, 1)) %>%
  select(-state_value, -state_landings) %>%
  rename(State = state) %>%
  mutate(State = str_to_title(State))
```

Now, use `kable()` + `kableExtra` to nicely format it for HTML:

```
kable(state_table) %>%
  kable_styling(bootstrap_options = "striped",
                full_width = FALSE) %>%
  add_header_above(c("", "2015 US commercial fisheries - top 5 states by value" = 2))
```

## 10.12 Sync with GitHub remote

Stage, commit, (pull), and push your updates to GitHub for safe storage & sharing. Check to make sure that the changes have been stored.

### 10.12.1 End synthesis session!