

PSYC 259: Principles of Data Science

Week 2: Data Workflow

Today

1. Workflow examples
2. “Minimizing mistakes” article discussion
3. Project structure principles
4. File organization + Import tutorial
5. How to get help
6. Homework 1

Workflow examples

Example 1

It was the year 2006...

- A young lab manager was tasked with testing how people walking on a treadmill perceived their own walking speed compared to life-sized avatars walking alongside them on a projection screen
- Software was written
- Results were analyzed
- A poster was presented
- The files sat dormant in storage for 15 years, never seeing the light of day....until now!

BigWalkerSpeed1.doc	Exp1data	6191.txt	Participant 6192
BigWalkerSpeed08-16-07.doc	Exp1documents	6192.txt	Age 21
SpatialWalkerDynamics	experiment2	6193.txt	Sex male
WalkingPerception	experiment3	6194.txt	Order 2
	experiment4	6195.txt	First Speed 2.4
	experiment5position	output	6/19/2006 12:00:26 PM
	perspective	perspective.sav	*****
			Block 1
			2.4 mid
		1 slo slower True	
		2 slo slower True	
		3 fas faster True	
		4 slo slower True	
		5 slo slower True	
		6 slo faster False	
		7 slo faster False	
		8 slo slower True	
		9 fas faster True	
		10 fas faster True	
		11 fas faster True	
		12 fas faster True	
		13 fas faster True	
		14 slo slower True	
		15 slo slower True	
		16 fas slower False	
		17 fas faster True	
		18 fas faster True	
			6192.txt
			Plain Text Document - 3 KB
			Information
		Created	Monday, June 19, 2006 at 10:10 AM
		Modified	Monday, June 19, 2006 at 10:10 AM

▶ Exp1data
▶ Exp1documents
▶ experiment2
▶ experiment3
▶ experiment4
▶ experiment5position
▶ perspective
▶ 6191.txt
▶ 6192.txt
▶ 6193.txt
▶ 6194.txt
▶ 6195.txt
▶ output
▶ perspective.sav

▶ 6192.txt
▶ 6193.txt
▶ 6194.txt
▶ 6195.txt
▶ 6211.txt
▶ 6221.txt
▶ 6222.txt
▶ 6223.txt
▶ 6224.txt
▶ 6231.txt
▶ 6232.txt

Block 1
2.4 mid

1	slo	slower	True
2	slo	slower	True
3	fas	faster	True
4	slo	slower	True
5	slo	slower	True
6	slo	faster	False
7	slo	faster	False
8	slo	slower	True
9	fas	faster	True
10	fas	faster	True
11	fas	faster	True
12	fas	faster	True
13	fas	faster	True
14	slo	slower	True
15	slo	slower	True
16	fas	slower	False
17	fas	faster	True
18	fas	faster	True
19	slo	slower	True
20	fas	faster	True

Correct: 17
Slower: 9

Block 2
2.7 near

1	slo	slower	True
2	slo	slower	True
3	slo	slower	True
4	slo	slower	True

id	age	sex	order	firstblock	near24	mid24	far24	near27	mid27	far27	near24s	mid24s	far24s	near27s	mid27s	far27s
6191	25	1	1	1	16	15	16	18	17	15	10	7	12	10	11	13
6192	21	1	2	1	18	17	18	18	19	17	8	9	12	8	11	9
6193	21	1	3	1	17	12	18	18	17	19	7	6	10	8	9	9
6194	22	1	1	2	17	17	17	20	19	17	7	7	11	10	9	11
6195	19	2	2	2	15	17	14	17	17	18	7	9	16	7	7	12
6211	19	1	3	2	20	16	13	17	16	18	10	14	13	13	8	10
6221	19	2	1	1	16	12	18	15	16	12	8	6	12	7	10	10
6222	21	2	2	1	20	13	18	19	18	18	10	13	12	9	8	12
6223	18	2	3	1	17	15	15	18	17	17	7	7	15	8	9	9
6224	18	2	1	2	19	15	17	16	19	16	11	9	13	10	11	8
6231	22	2	2	2	15	15	16	20	16	19	9	9	10	10	8	11
6232	21	1	3	2	14	14	11	14	16	15	4	8	7	6	8	13

a id	age	sex	order	firstblock	near24	mid24	far24	near27	mid27	far27	near24s	mid24s	far24s	near27s	mid27s	far27s
6191	25	1	1	1	16	15	16	18	17	15	10	7	12	10	11	13
6192	21	1	2	1	18	17	18	18	19	17	8	9	12	8	11	9
6193	21	1	3	1	17	16	18	18	17	19	7	6	10	8	9	9
6194	22	1	1	2	17	17	17	20	19	17	7	7	11	10	9	11
6195	19	2	2	2	15	17	14	17	17	18	7	9	16	7	7	12
6211	19	1	3	2	20	16	13	17	16	18	10	14	13	13	8	10
6221	19	2	1	1	16	12	18	15	16	12	8	6	12	7	10	10
6222	21	2	2	1	20	13	18	19	18	18	10	13	12	9	8	12
6223	18	2	3	1	17	15	15	18	17	17	7	7	15	8	9	9
6224	18	2	1	2	19	15	17	16	19	16	11	9	13	10	11	8
6231	22	2	2	2	15	15	16	20	16	19	9	9	10	10	8	11
6232	21	1	3	2	14	14	11	14	16	15	4	8	7	6	8	13

Block 1 2.4 mid				
1	slo	slower	True	
2	slo	slower	True	
3	fas	faster	True	
4	slo	slower	True	
5	slo	slower	True	
6	slo	faster	False	
7	slo	faster	False	
8	slo	slower	True	
9	fas	faster	True	
10	fas	faster	True	
11	fas	faster	True	
12	fas	faster	True	
13	fas	faster	True	
14	slo	slower	True	
15	slo	slower	True	
16	fas	slower	False	
17	fas	faster	True	
18	fas	faster	True	
19	slo	slower	True	
20	fas	faster	True	
Correct:	17			
Slower:	9			

per27n	per27m	per27f	mean24	mean27	data	plspeed	plc ond	
90.00	85.00	75.00	78.33	83.33	80.00	2.40	near	
90.00	95.00	85.00	88.33	90.00	90.00	2.40	near	
90.00	85.00	95.00	78.33	90.00	85.00	2.40	near	
100.00	95.00	85.00	85.00	93.33	85.00	2.40	near	
85.00	85.00	90.00	76.67	86.67	75.00	2.40	near	
85.00	80.00	90.00	81.67	85.00	100.00	2.40	near	
75.00	80.00	60.00	76.67	71.67	80.00	2.40	near	
95.00	90.00	90.00	85.00	91.67	100.00	2.40	near	
90.00	85.00	85.00	78.33	86.67	85.00	2.40	near	
80.00	95.00	80.00	85.00	85.00	95.00	2.40	near	
100.00	80.00	95.00	76.67	91.67	75.00	2.40	near	
70.00	80.00	75.00	65.00	75.00	70.00	2.40	near	
-	-	-	-	-	75.00	2.40	mid	
-	-	-	-	-	85.00	2.40	mid	
-	-	-	-	-	60.00	2.40	mid	
-	-	-	-	-	85.00	2.40	mid	
-	-	-	-	-	85.00	2.40	mid	
-	-	-	-	-	80.00	2.40	mid	
-	-	-	-	-	60.00	2.40	mid	
-	-	-	-	-	65.00	2.40	mid	
-	-	-	-	-	75.00	2.40	mid	
-	-	-	-	-	75.00	2.40	mid	
-	-	-	-	-	75.00	2.40	mid	
-	-	-	-	-	70.00	2.40	mid	
-	-	-	-	-	80.00	2.40	far	
-	-	-	-	-	90.00	2.40	far	
-	-	-	-	-	90.00	2.40	far	
-	-	-	-	-	85.00	2.40	far	
					70.00	2.40	far	

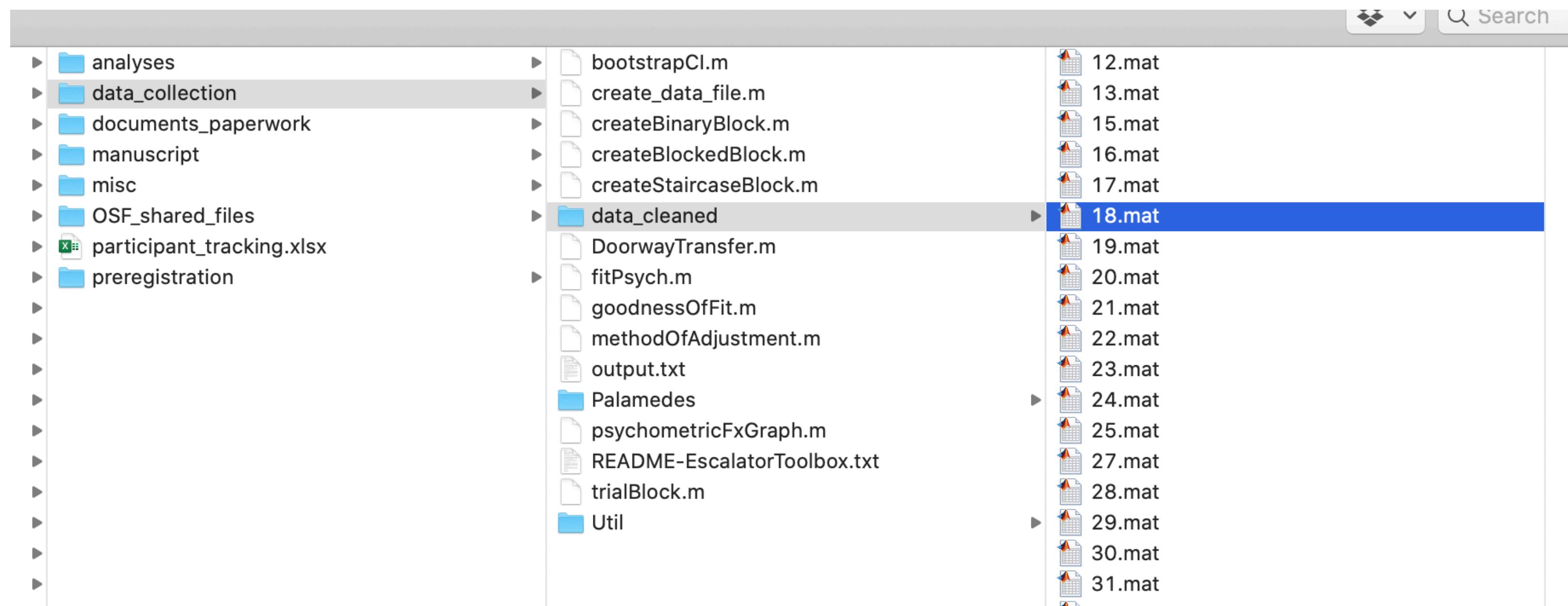
List of crimes against data science

- Files: unorganized, inconsistently named, used unsupported proprietary formats, no version control, cluttered workspace with intermediary/redundant file
- Automation (lack of)
 - Used copy/paste to transfer data rather than automating input
 - Copied data within SPSS from wide to long
 - Made graphs/analyses from drop-down menus
 - No metadata, no documentation

Example 2

The year was 2019...

- An assistant professor had a team of 5 undergrad RAs run participants in a study on perception of fitting through doorways
- The RAs entered the participant and condition info into a Matlab program, which ran the trials and wrote output data for each participant
- The RAs logged metadata into a “participant_tracking.xls” document after running each participant, leaving notes about problems with each participant



Command Window

```

>> sex

sex =
    1x1 cell array

    {'m'}

>> id

id =
    1x1 cell array

    {'18'}

>> block1

block1 =
    struct with fields:

        mode: 'staircase'
        start_unit: 32
        down_step: 4
        up_step: 3
        num_trials: 16
        use_for_fit: 1

>> output_aff_f

output_aff_f =
    struct with fields:

        id: '18'
        condition: 'aff_f'
        threshold: 40.8259
        slope: 1.0576
        trial_unit: [1x30 double]
        trial_unit_fit: [1x30 double]
        trial_resp: [1x30 double]
        trial_resp_fit: [1x30 double]
        trial_mode: [1x30 double]
        trial_subblock: [1x30 double]
        stim_levels: [1x107 double]
        elapsed_time: 618.8013

```

Workspace

Name	Value
block1	1x1 struct
block2	1x1 struct
block_fitting	2x1 cell
block_squeezing	2x1 cell
door_measure1	1x1 cell
door_measure2	1x1 cell
height1	1x1 cell
height2	1x1 cell
id	1x1 cell
judgment_cond	1x1 cell
measurement_...	1x1 cell
out_dir	'C:\Users\labuse...
output_aff_f	1x1 struct
output_aff_s	1x1 struct
output_pre_moa	1x1 struct
output_pst_moa	1x1 struct
practice_cond	1x1 cell
<input checked="" type="checkbox"/> save_figs	1
save_file	'C:\Users\labuse...
sex	1x1 cell
stim_levels	1x107 double
weight1	1x1 cell
weight2	1x1 cell

```

id,sex,judgecond,practcond,height,weight,doormeasure,aff_s,span_s,aff_f,span_f,pre1,pre2,pre3,pre4,pre5,pre6,pst1,pst2,
pst3,pst4,pst5,pst6
100,f,S,S,60,415,47,3,25,3,19,-0,5,41,8,40,15,44,2,43,1,44,35,41,9,43,5,36,05,31,6,30,55,27,75,30,25,25
101,m,F,F,67,65,64,7,24,7,27,6,1,45,1,3,49,5,46,5,40,45,44,2,43,35,43,05,48,45,47,9,47,25,49,3,47,75,6
102,f,S,S,62,55,42,9,22,55,25,1,0,5,41,1,12,44,95,41,2,42,42,2,42,75,42,55,25,45,25,7,23,35,24,4,24,25,23,85
103,m,F,F,71,68,28,15,29,4,1,44,9,45,95,46,85,44,7,46,6,44,15,43,6,47,1,44,7,44,65,44,5,43,65,45,1
104,m,S,F,70,117,6,34,9,36,6,2,54,6,6,38,65,37,85,41,85,42,55,37,8,43,7,38,65,35,45,37,55,38,35,37,75,37,8
105,f,S,S,57,9,42,9,23,95,24,2,3,5,45,2,5,35,65,32,8,33,45,35,25,37,9,33,8,25,65,24,7,27,45,27,5,23,6,25,7
106,f,F,S,65,5,81,7,31,6,31,1,0,5,47,6,7,50,7,51,35,55,2,51,25,52,8,48,4,48,7,48,05,51,8,48,15,49,6
107,f,F,F,62,05,48,1,25,4,30,7,1,5,42,2,1,44,8,44,65,44,75,43,85,42,95,42,3,43,95,42,5,43,75,44,65,46,4,46,2
108,f,S,S,71,1,120,6,37,35,28,1,1,5,57,6,6,5,54,3,53,75,51,15,49,6,52,2,54,35,57,8,56,45,53,6,56,3,53,05,57,3
109,m,F,F,70,6,75,5,39,2,31,0,45,5,1,5,34,95,37,3,38,1,35,95,38,3,35,7,40,35,42,5,43,9,39,15,41,8,42,75
nboot = 112,F,S,F,71,75,110,8,41,05,37,5,5,61,5,2,46,4,43,65,48,35,43,75,40,8,42,2,44,7,49,5,39,9,40,55,43,55,37,9
113,M,F,F,71,79,27,65,24,5,0,45,4,50,6,48,75,58,7,49,1,49,4,49,35,51,7,48,7,49,05,48,2,48,4,48,5
for i =
if
115,m,S,F,73,35,79,25,27,15,33,1,0,5,43,6,2,36,5,35,35,1,44,44,3,43,6,45,4,39,55,41,45,40,15,40,6,42,1
116,m,S,F,73,94,7,27,1,25,5,2,5,46,5,5,32,15,28,85,32,55,31,6,29,9,29,9,36,15,32,4,30,25,32,05,31,9,30,95
12,f,F,S,61,54,2,25,15,28,6,0,5,43,1,4,5,50,54,4,52,5,53,55,53,9,54,55,45,5,47,2,45,8,45,85,47,8,47,45
13,f,F,S,62,5,49,6,25,85,23,8,1,5,39,8,3,5,44,1,44,4,44,45,2,42,7,43,65,37,05,39,2,40,15,40,4,38,8,38,65
15,m,F,F,67,95,55,1,26,05,27,8,0,45,8,5,45,47,05,49,8,51,5,55,51,8,47,7,47,7,51,55,50,2,51,9,51,5
16,m,S,F,NaN,NaN,NaN,33,2,3,5,47,7,9,5,35,95,36,75,33,55,37,25,36,4,39,95,43,9,46,15,46,43,9,45,75,45,4
17,f,F,S,66,9,48,6,36,2,28,2,3,39,7,3,42,35,43,35,41,7,39,35,34,9,36,95,48,47,7,49,35,49,6,47,95,47
18,m,S,S,67,6,63,9,38,15,33,7,9,42,2,2,53,45,53,2,58,4,54,9,53,8,51,7,48,7,40,1,39,75,40,4,39,25
19,f,F,F,60,4,84,5,34,05,36,1,1,5,56,6,4,56,75,52,15,53,9,55,15,55,95,57,95,49,2,46,85,47,1,49,15,50,7,49,85
20,m,S,F,67,5,86,4,27,9,25,3,2,5,47,8,5,43,95,44,55,40,55,41,75,40,5,40,1,34,75,36,7,34,05,35,5,33,6,33,45
21,f,F,S,59,25,45,8,25,75,26,6,1,5,39,1,6,41,55,38,6,39,2,38,9,41,15,40,7,42,8,38,4,39,35,39,85,38,7,40,7
22,m,S,S,66,5,75,7,28,05,30,6,-0,5,44,6,1,38,75,38,45,39,7,36,2,38,3,39,8,40,1,40,35,41,9,44,45,43,5,44,3
23,f,F,F,65,45,76,8,32,55,29,6,4,53,6,4,5,45,85,47,9,50,2,49,75,48,75,51,55,54,35,53,15,58,8,51,95,52,65,53,75
24,f,S,F,62,75,52,9,22,95,25,1,43,3,5,45,95,39,38,85,37,75,37,8,36,55,36,35,35,9,35,3,36,25,36,15,33
25,m,F,S,66,2,91,8,42,7,37,4,49,5,-0,5,57,65,59,2,56,75,58,2,59,45,63,5,58,75,58,8,58,9,58,8,57,85
27,m,F,F,71,5,66,7,36,45,28,1,43,10,49,3,54,5,58,1,51,45,53,3,56,85,54,25,54,95,52,95,55,35,54,05,56,95
28,m,S,S,73,4,74,6,25,9,26,9,1,42,4,5,5,33,65,33,1,32,95,34,05,36,55,39,6,29,25,30,15,28,55,29,85,27,9,27,15
29,f,S,S,61,3,59,2,27,55,25,6,0,42,1,6,43,5,43,05,43,5,45,9,47,95,47,5,28,75,29,45,32,7,38,5,32,1,28,15
30,f,F,F,60,85,62,25,85,28,7,0,5,45,7,12,5,48,35,46,35,47,75,49,45,49,2,51,15,47,9,50,44,25,49,95,48,2,48,5
31,m,S,F,71,6,76,1,25,95,28,1,1,5,42,6,4,41,15,42,85,45,75,42,2,43,7,39,3,36,25,39,42,25,38,7,39,9,40,25
32,m,F,S,71,6,77,27,65,25,6,0,42,1,3,5,51,65,49,35,49,2,46,65,45,1,43,5,47,05,47,5,44,45,46,55,45,15,45,55
33,m,S,S,72,5,83,3,27,15,29,3,6,42,8,3,46,35,43,05,41,1,41,37,55,40,5,30,95,31,1,31,2,29,8,31,25,28,25

```

- “create_data_file.m” loads each individual file in “data_cleaned” and creates rows of data to write to “output.txt”
- “output.txt” gets copied to “analyses” for R analyses

The screenshot shows the RStudio interface with the following components:

- Top Bar:** Shows the file "analyses.R" is open. Buttons for "Run", "Source", and "History" are visible.
- Left Panel (Code View):** Displays the R script "analyses.R". The code reads a CSV file "output.txt", creates factor levels for "id", "practice", "judge", and "match", and calculates error terms for each judge category. It also defines variables for pre-error terms and past-error terms.
- Global Environment:** Shows objects "ds" (100 obs. of 26 variables), "my_theme" (List of 11), and "pd" (Environment). A variable "cbp1" is highlighted.
- File Browser:** Shows the project structure under "analyses" folder in "Dropbox/past_projects/study_doorwaytransfer". Files listed include ".RData", ".Rhistory", "age.xlsx", "analyses.R", "analyses.Rproj", "correlations.pdf", "error.pdf", "output.txt", and "summarySE.R".

```

30
31 # READ DATA ----
32
33 #Read data and set factor levels
34 ds <- read_csv("output.txt")
35 ds$id <- factor(ds$id)
36 ds$practice <- factor(ds$practcond, levels =c("S","F"),labels = c("Squeezing","Fitting"))
37 ds$judge <- factor(ds$judgecond, levels =c("S","F"),labels = c("Squeezing","Fitting"))
38 ds$match <- ds$practcond == ds$judgecond
39 ds$match <- factor(ds$match, levels = c(T, F), labels = c("Congruent","Incongruent"))
40 |
41 # Create error terms ====
42 ds$judge_ref[ds$judge == "Fitting"] <- ds$aff_f[ds$judge == "Fitting"]
43 ds$judge_ref[ds$judge == "Squeezing"] <- ds$aff_s[ds$judge == "Squeezing"]
44
45 ds$pre_err1 <- ds$pre1 - ds$judge_ref
46 ds$pre_err2 <- ds$pre2 - ds$judge_ref
47 ds$pre_err3 <- ds$pre3 - ds$judge_ref
48 ds$pre_err4 <- ds$pre4 - ds$judge_ref
49 ds$pre_err5 <- ds$pre5 - ds$judge_ref
50 ds$pre_err6 <- ds$pre6 - ds$judge_ref
51 ds$pst_err1 <- ds$pst1 - ds$judge_ref
52 ds$pst_err2 <- ds$pst2 - ds$judge_ref
53 ds$pst_err3 <- ds$pst3 - ds$judge_ref
54 ds$pst_err4 <- ds$pst4 - ds$judge_ref
55
40:1 # READ DATA

```

- R reads in output.txt, assigns factors, and calculates DVs
- Creates the analyses and figures for the paper

doorway_transfer.tex

ION1.pdf

apparatus-compressed.pdf
apparatus.pdf
correlations.pdf
DesignFigure.pdf
diff.aux
diff.bbl
diff.blg
diff.log
diff.out
diff.pdf
diff.synctex.gz
diff.tex
doorway_transfer.aux
doorway_transfer.bbl
doorway_transfer.blg
doorway_transfer.log
doorway_transfer.out
doorway_transfer.pdf
doorway_transfer.synctex.gz
doorway_transfer.tex
error.pdf
ExampleTrials.mp4
master.bib
sage_latex_template_4
SageH bst
sagej.cls
sagej.log
SageV bst
submission1
submission2
submission3-publication

Markup More...

Typeset LaTeX Macros Tags Labels Templates

replicability \citep{AsendorptConner2013, SimmonsNelson2011}, the current study's procedure and analysis plan were pre-registered before data collection began. The pre-registration document was entered on AsPredicted.com (\url{https://aspredicted.org/s58hb.pdf}). The sample size was determined in advance based on a power analysis: Past work comparing pretest/posttest affordance judgment errors across multiple between-subjects conditions found a large interaction effect size, $\eta^2 = .152$ \citep{Recal}. However, the transfer effect in the current study might be smaller, so the sample size was conservative. The resulting effect size ($\eta^2 = .02$) resulted from a technical error in the pre-registration. Although postural sway measurements from the accelerometer data were measured relative to trials and not absolute, the pre-registration we recorded the absolute error. In this paper, we refer to this as "mismatching"; here we report the absolute error.

error.pdf

PDF document - 5 KB

Information Show More

Created Aug 7, 2019 at 3:42 PM
Modified Aug 7, 2019 at 3:42 PM

Tags Add Tags...

Absolute error (cm)

Pretest Phase Posttest

Practiced Action
Congruent (Orange circle)
Incongruent (Blue circle)
Judged Action
Squeezing (Orange square)
Fitting (Blue square)

12 Quarterly Journal of Experimental Psychology XX(X)

104
105 \subsection{Participants and Materials}

106
107 The final sample included 54 participants (46 female, 8 male) assigned to one of two conditions in a 2 Judged Action (congruent, incongruent) \times 2 Practiced Action (squeezing-congruent, squeezing-incongruent) \times 2 Judged Phase (male, female) design. There were 11 participants in each of the four conditions. The final sample included 54 participants (46 female, 8 male) assigned to one of two conditions in a 2 Judged Action (congruent, incongruent) \times 2 Practiced Action (squeezing-congruent, squeezing-incongruent) \times 2 Judged Phase (male, female) design. There were 11 participants in each of the four conditions. Judgment phase was a within-subjects factor. Two additional participants were excluded for failure to follow instructions and one participant was replaced for failure to follow repeated requests from the experimenter.

Absolute error (cm)

Pretest Posttest

Practiced Action
Congruent (Orange circle)
Incongruent (Blue circle)
Judged Action
Squeezing (Orange square)
Fitting (Blue square)

Figure 3. Absolute judgment errors by phase (x-axis: pretest vs posttest), judged action (circles: squeezing; squares: fitting), and practiced action (orange symbols: congruent practice; blue symbols: incongruent practice).

with the pre-registration plan, participants with outlier data (based on inter-quartile range) were excluded and not replaced (3 participants in the FC condition and 1 participant in the FI condition). The influence of outliers on the results is discussed below.

Overall model Table 1 shows results for the overall model predicting absolute judgment errors from judged action, practiced action, and judgment phase (and their interactions) as fixed effects and random intercepts by participant (random slope models failed to converge). As seen in Figure 3, a significant main effect of phase indicates that participants were more accurate in posttest compared with pretest, and a significant main effect of judged action reveals that participants were more accurate when judging fitting compared with squeezing. However, these effects were moderated by significant two-way (judged action \times phase, practiced action \times phase) and three-way (judged action \times practiced action \times phase) interactions. Including or excluding outliers did not affect the significance of any of the effects in the model.

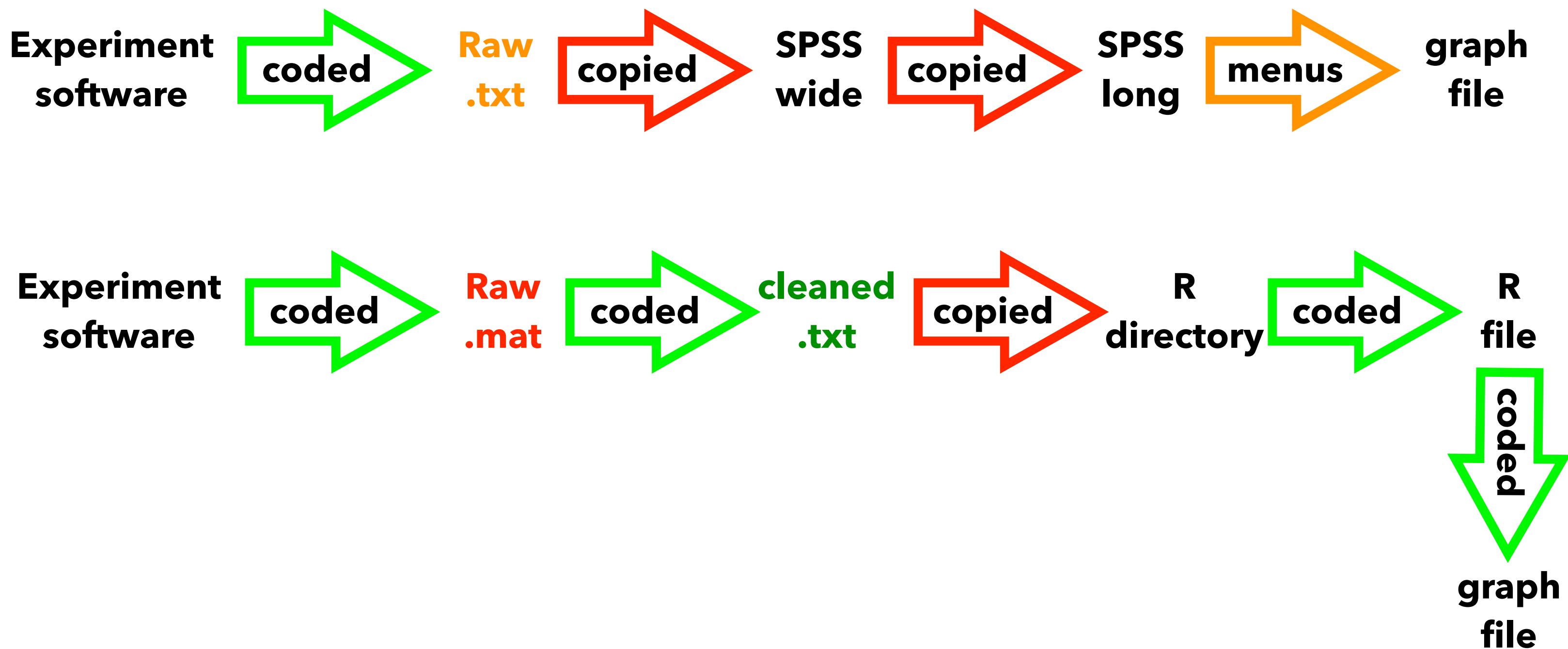
What was good?

- Files: Organized in directories, consistently named, mostly flat text files
- Version control: minimal intermediaries, used GitHub
- Automation: Data collection was automated, cleaning raw data was automated, analyses/figures were automated
- Automation: Minimal copy/paste, minimal user typing
- Metadata saved about each participant to diagnose problems, understand exclusion

What should be better?

- Using matlab scripts/data files makes project harder to share, less future-proof; no way to archive the coding environment
- Copy/pasting data/figures between folders, analyses to paper leave plenty of room for error
- Clear set of processes to regenerate from raw data, but not documented well (requires user to run scripts in correct order which only I know)
- No formal checks of data quality
- Lots of repetitive, single-use Matlab and R code

Comparing the workflows



Minimizing mistakes
article discussion

What principles should
guide project structure?

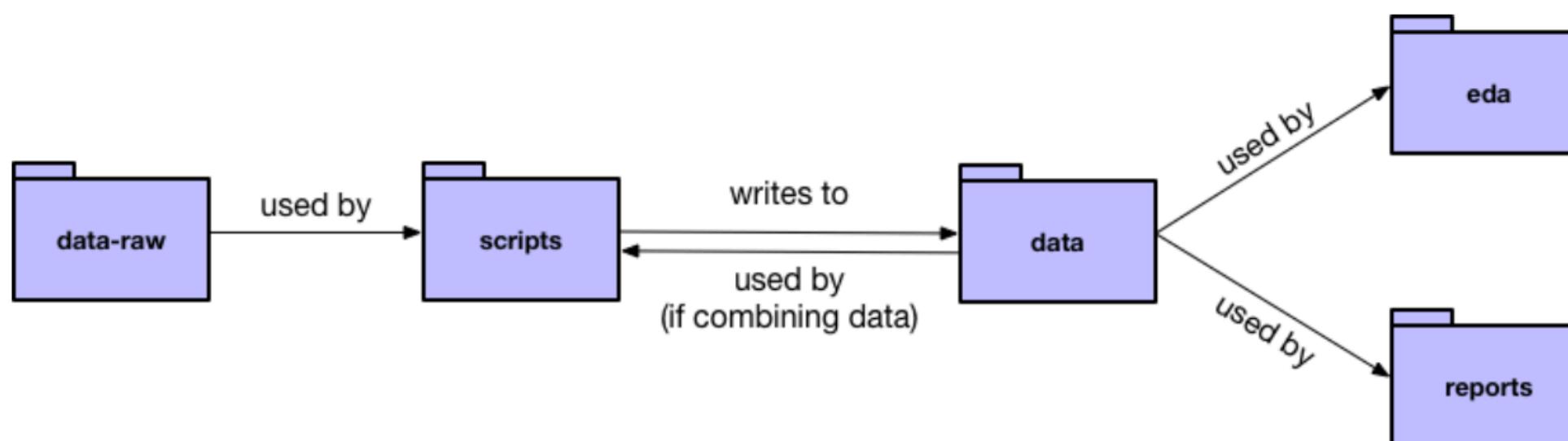
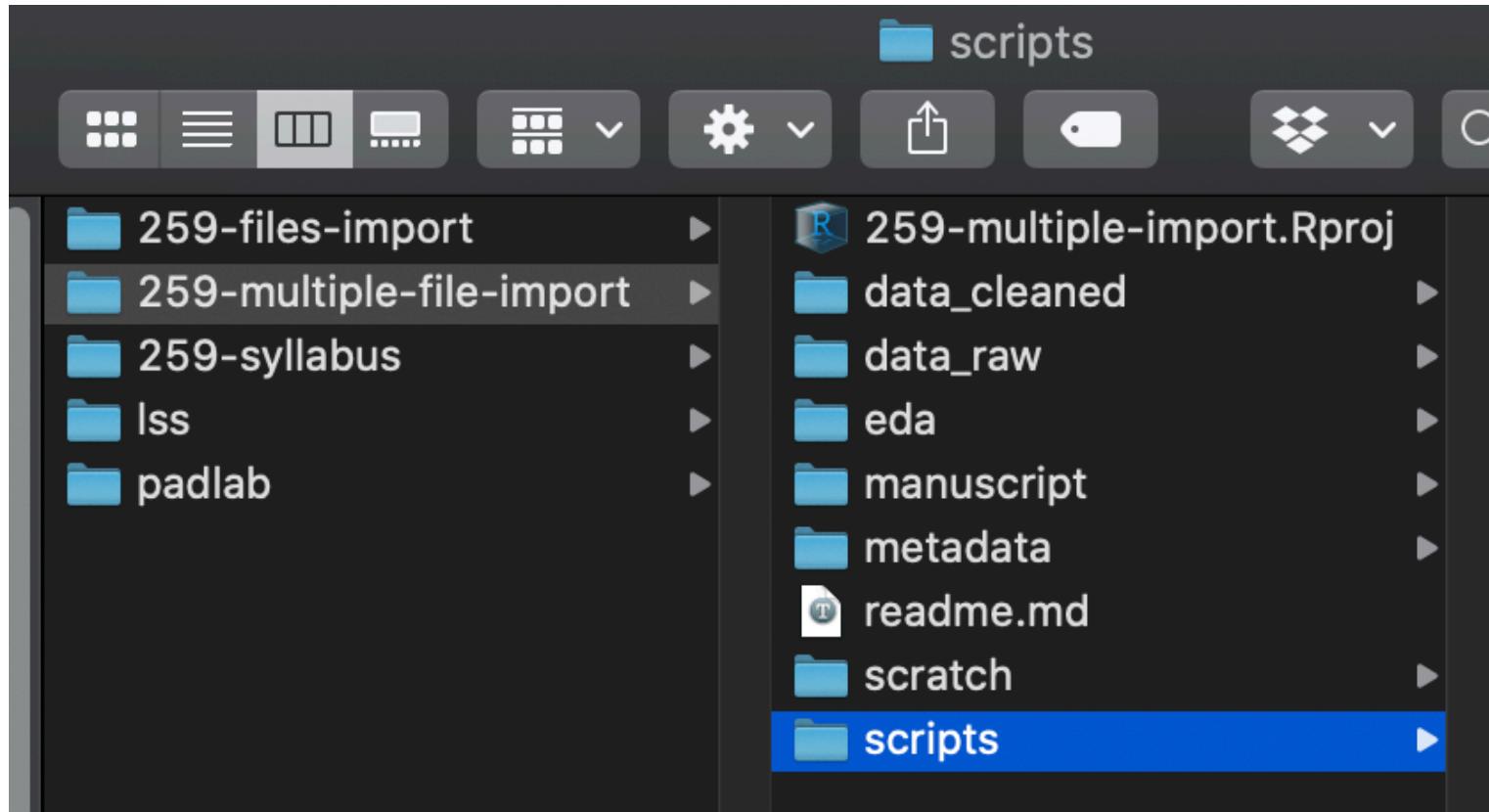
Two principles of project workflow

1. Folder organization creates rules and defines a workflow; establishes a location from which to build relative file paths
2. Version control tracks file history without duplication/clutter; allows for collaboration/derivation/experimentation

What principles should
guide project structure?

#1 File/folder organization

Folders to organize similar file types (w/in a root project folder)



R Studio Project

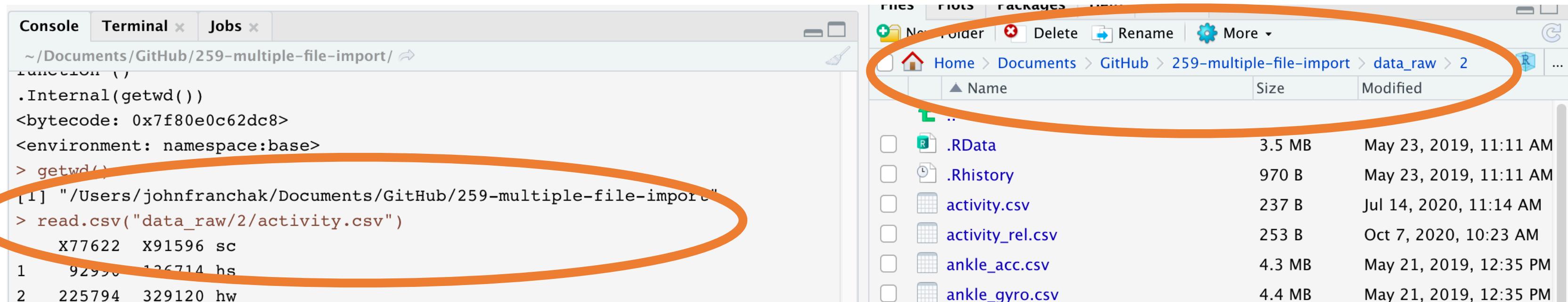
- Essentially, your project folder with some files that keep tabs on your command history
 - Can create a project from an existing folder
 - Open the RProj file to open the project
 - Can save your workspace, but follow R4DS guidelines and decline that feature
- What is real?
 - Not real: objects/data frames in your workspace
 - Real: data files and the scripts used to work with them
 - Intermediaries incur a cost of maintenance

Accessing files from your working directory

- Your RStudio Project folder is your default *working directory*
 - working directory: where R can find files
 - `getwd()` probably gives you something ugly like `/Users/johnfranchak/Documents/GitHub/259-multiple-file-import`
- *Absolute file paths* like that should be avoided at all costs!
 - Absolute file paths don't transfer between computers with different user names or different operating systems
 - Bad for reproducibility, extensibility, and sharing
- But don't you need them to get into all of those directories you just told me I need to use?

Relative files paths 💪💪💪

- Just tell R what to look for relative to your project (working) directory!



```
Console Terminal x Jobs x
~/Documents/GitHub/259-multiple-file-import/ ↵
[1] .Internal(getwd())
<bytecode: 0x7f80e0c62dc8>
<environment: namespace:base>
> getwd()
[1] "/Users/johnfranchak/Documents/GitHub/259-multiple-file-import"
> read.csv("data_raw/2/activity.csv")
   X77622  X91596 sc
1  92950  126714 hs
2  225794  329120 hw
```

Name	Size	Modified
.RData	3.5 MB	May 23, 2019, 11:11 AM
.Rhistory	970 B	May 23, 2019, 11:11 AM
activity.csv	237 B	Jul 14, 2020, 11:14 AM
activity_rel.csv	253 B	Oct 7, 2020, 10:23 AM
ankle_acc.csv	4.3 MB	May 21, 2019, 12:35 PM
ankle_gyro.csv	4.4 MB	May 21, 2019, 12:35 PM

- *here* package detects the project directory and composes filenames from/to any folder
 - `here("folder1", "folder2", "filename")`

Other considerations

- Not everything can be automated, but you can try to limit human data entry to a single master table
 - Track notes about sessions, inclusion/exclusion info
 - Keep as part of project metadata, and use it to direct your scripts to pull the 'right' data
- Not every project can be contained in a local directory on a single computer (or on github)
 - Large datasets might need other solutions
 - "data_raw" might need to be "data_slightly_cooked"

Tutorial #1 - Data Import/Export

- R Studio Cloud/Github repos
 - 259-files-import
 - 259-multiple-file-import

Assumptions

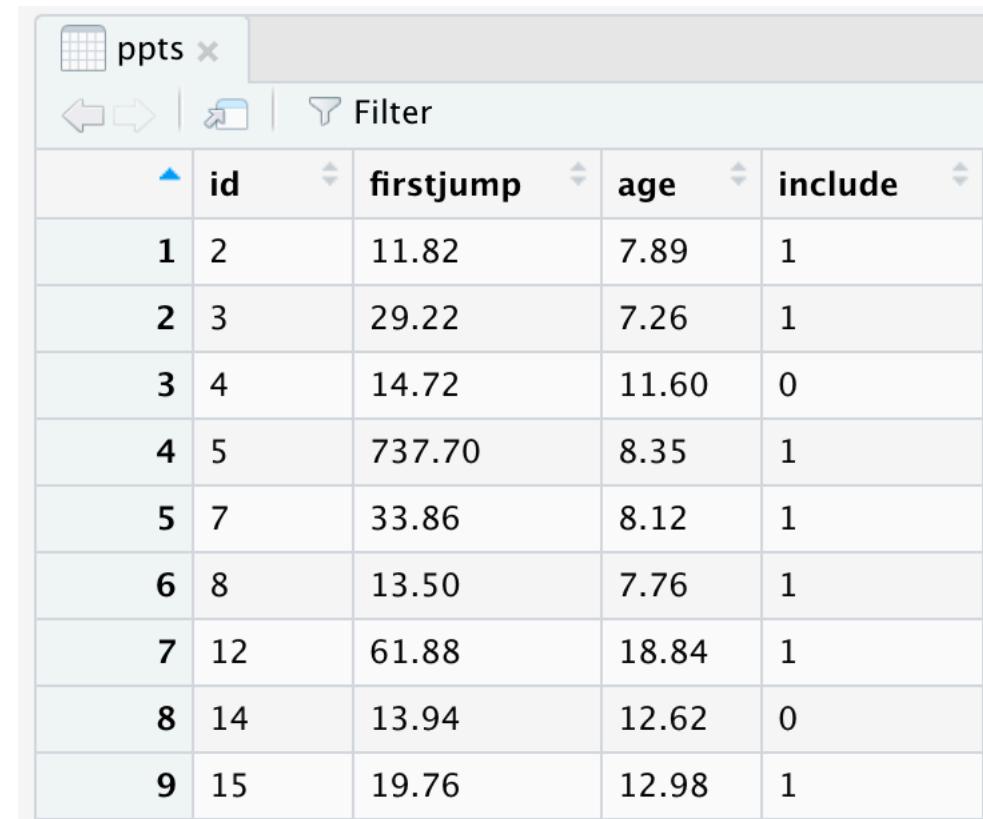
- Loading packages
- Running chunks of code from a script
- Basic syntax for assigning variables <-
- Objects/data get stored in the environment
- What a dataframe/tibble is

readr package (loads via “tidyverse”)

- `read_csv()`, `read_tsv()`, `read_delim()`
 - Variants to read flat file data w/ different delimiters
 - Reads to a tibble rather than a data frame
 - Doesn't automatically convert strings to factors
 - Faster than their base r counterparts (e.g., `read.csv`)
- `vroom`
 - Reads multiple files to a single dataset
 - Guesses the delimiter

What's a tibble?

- A type of data frame: rectangular data with columns (variables of different types) and rows (usually observations)
- A data frame is a spreadsheet, but you can't edit it!



A screenshot of a data editor window titled "ppts". The window shows a table with 9 rows and 5 columns. The columns are labeled "id", "firstjump", "age", and "include". The data is as follows:

	id	firstjump	age	include
1	2	11.82	7.89	1
2	3	29.22	7.26	1
3	4	14.72	11.60	0
4	5	737.70	8.35	1
5	7	33.86	8.12	1
6	8	13.50	7.76	1
7	12	61.88	18.84	1
8	14	13.94	12.62	0
9	15	19.76	12.98	1

What's a tibble?

- Tibbles are data frames that:
 - Don't do "partial matching"
 - Print a bit cleaner in your console
 - You will rarely notice the difference, and you can easily `as_tibble()` or `as.data.frame()` back and forth

```
> ppts_df <- as.data.frame(ppts)
> ppts_tb <- as_tibble(ppts)
> ppts_df$ag
[1]  7.89  7.26 11.60  8.35  8.12  7.76 18.84 12.62 12.98      NA 17.95 12.03 10.78 14.37 12.23 14.66 11.54
> ppts_tb$ag
NULL
Warning message:
Unknown or uninitialised column: `ag`.
```

Writing data

- `write_csv()`, `write_delim()`
 - Writes flat files from a data frame
- `save()`
 - Saves workspace variables as RData files
- `ggsave()`
 - Save a ggplot to an image file

RStudio Cloud

- Find the 259-files-import project and open it
 - It's an "assignment" (but not really) which means you can open it and it creates a copy that you can work in
 - I encourage you to open it up and follow along
 - Add comments/notes to the code (if that helps)
 - Technical problems -> chat with Jake and visit a break-out room

How to get help

Where to search for help

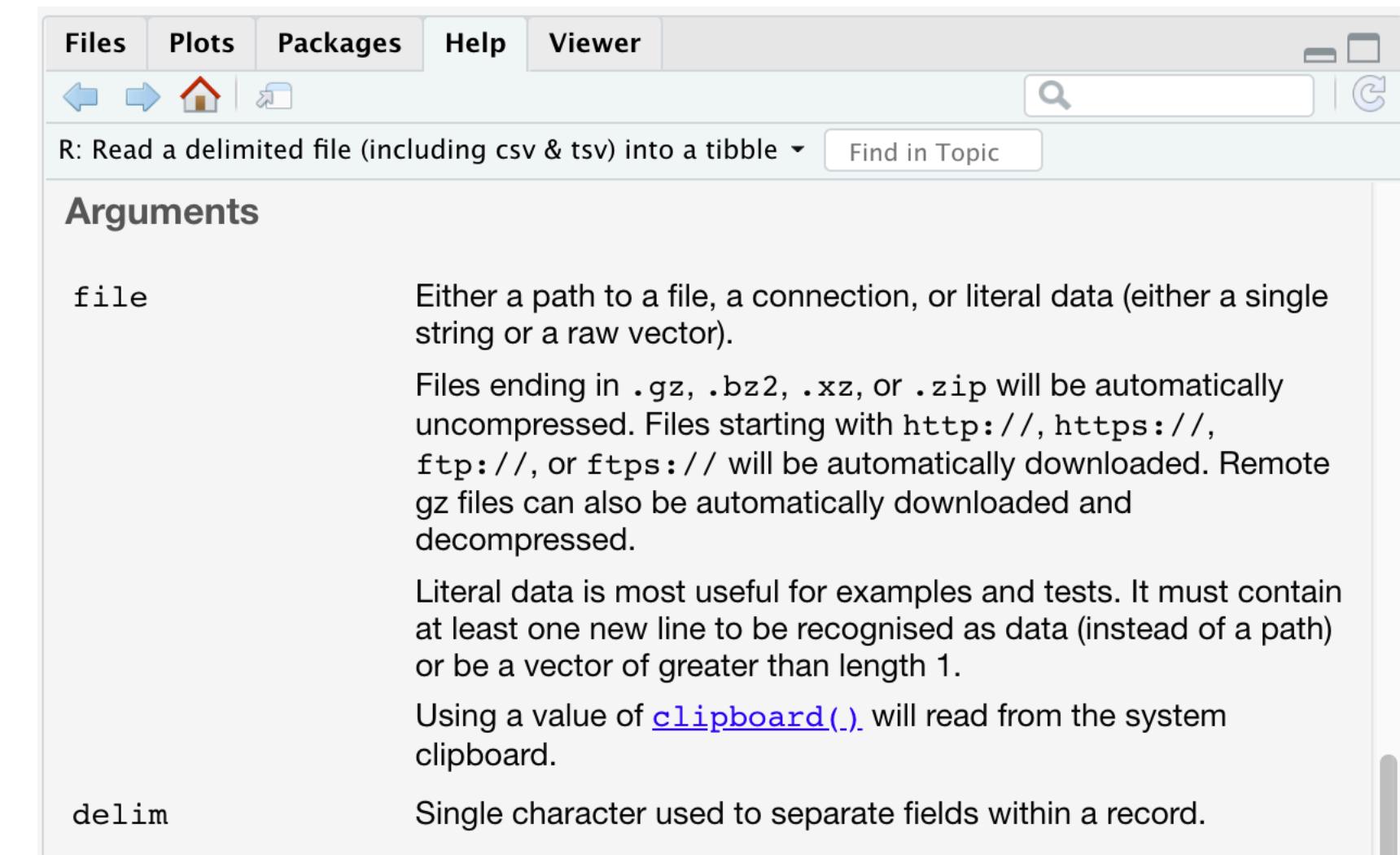
- R Documentation
 - ?function brings up the documentation for a function (assuming the package is loaded)

The screenshot shows the RStudio Help Viewer interface. The title bar includes tabs for Files, Plots, Packages, Help, and Viewer, with Help selected. Below the tabs are navigation icons for back, forward, home, and search, along with a magnifying glass icon for the search bar. The search bar contains the text "R: Read a delimited file (including csv & tsv) into a tibble". To the right of the search bar is a "Find in Topic" button. The main content area displays the documentation for the `read_delim` function. The title is `read_delim {readr}`. On the right, there is a link labeled "R Documentation". The documentation text starts with a brief description of the function's purpose: "Read a delimited file (including csv & tsv) into a tibble". It then provides a "Description" section explaining that `read_csv()` and `read_tsv()` are special cases of `read_delim()`, useful for reading CSV and TSV files respectively, and notes that the decimal separator is a semicolon (;) in some European countries. The "Usage" section shows the function signature:

```
read_delim(  
  file,  
  delim,  
  quote = "\"",  
  escape_backslash = FALSE,  
  escape_double = TRUE,  
  col_names = TRUE,  
  col_types = NULL,  
  locale = default_locale(),  
  na = c("", "NA"),  
  quoted_na = TRUE,  
  comment = "",  
  trim_ws = FALSE,  
  skip = 0,  
  n_max = Inf,  
  guess_max = min(1000, n_max),  
  progress = show_progress(),  
  skip_empty_rows = TRUE  
)
```

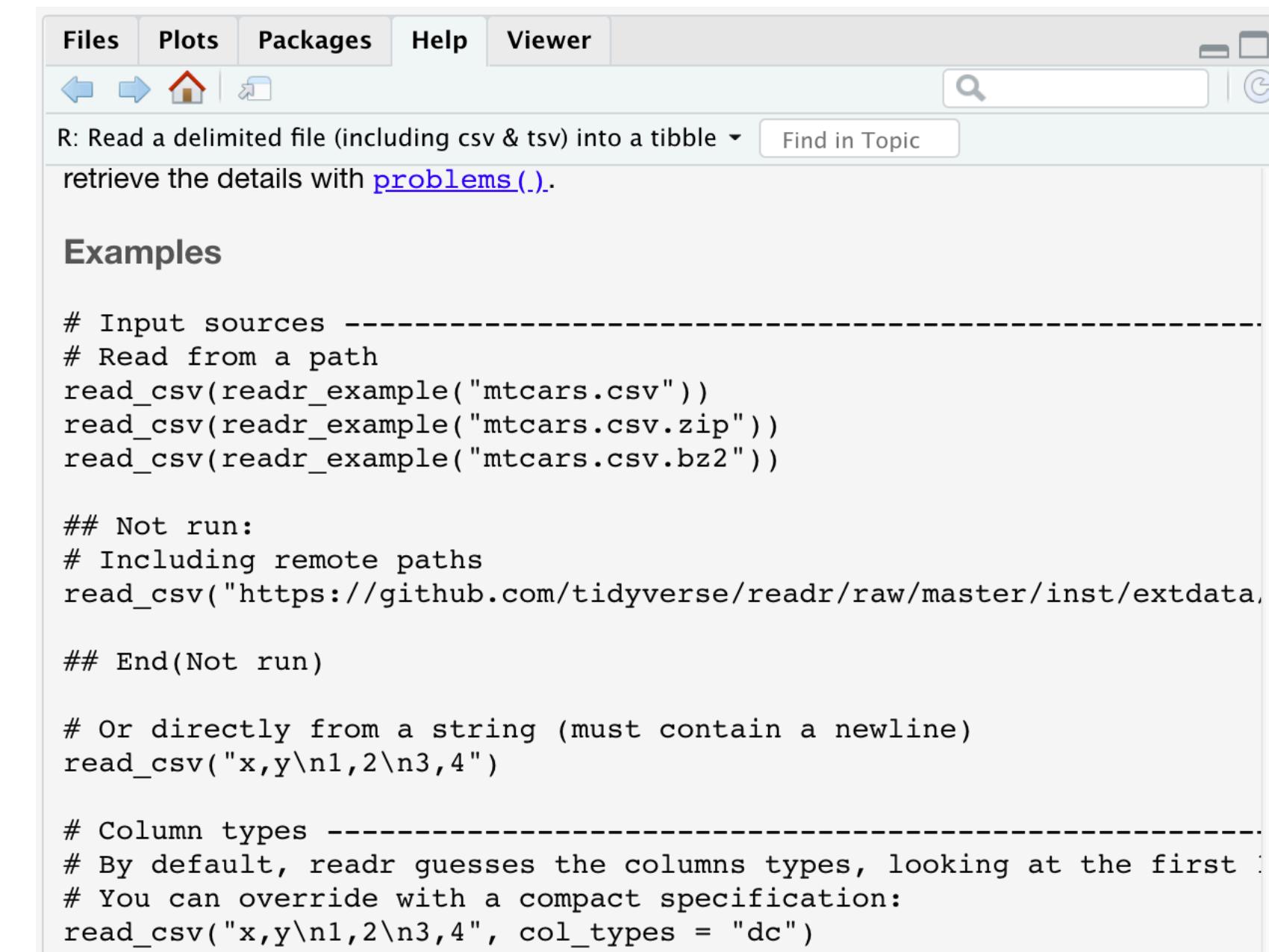
Where to search for help

- R Documentation
 - `?function` brings up the documentation for a function (assuming the package is loaded)
 - Scroll down to see the “arguments” definitions



Where to search for help

- R Documentation
 - ?function brings up the documentation for a function (assuming the package is loaded)
 - Scroll down to see the “arguments” definitions
 - Scroll down more to see examples



The screenshot shows the RStudio interface with the "Help" tab selected in the top menu bar. Below the menu, there are navigation icons for back, forward, and search, along with a search bar and a "Find in Topic" button. The main content area displays the documentation for the `read_csv` function. It starts with a brief description: "R: Read a delimited file (including csv & tsv) into a tibble". A "Find in Topic" button is also present here. Below the description, the word "retrieve" is highlighted with a blue underline. The "Examples" section contains several code snippets demonstrating how to use the function with different input sources and column types.

```
# Input sources -----
# Read from a path
read_csv(readr_example("mtcars.csv"))
read_csv(readr_example("mtcars.csv.zip"))
read_csv(readr_example("mtcars.csv.bz2"))

## Not run:
# Including remote paths
read_csv("https://github.com/tidyverse/readr/raw/master/inst/extdata/mtcars.csv")

## End(Not run)

# Or directly from a string (must contain a newline)
read_csv("x,y\n1,2\n3,4")

# Column types -----
# By default, readr guesses the columns types, looking at the first few lines
# You can override with a compact specification:
read_csv("x,y\n1,2\n3,4", col_types = "dc")
```

Where to search for help

- Package vignettes/blogs
 - <https://cran.r-project.org/web/packages/readr/vignettes/readr.html>
- Google
 - Blogs/instructional sites
 - StackExchange (especially for error text)
 - Tidyverse community
 - <https://community.rstudio.com/c/tidyverse/6>
- Take working examples and fiddle!

Tidyverse cheatsheets!

Data Import :: CHEAT SHEET

R's **tidyverse** is built around **tidy data** stored in **tibbles**, which are enhanced data frames.



The front side of this sheet shows how to read text files into R with **readr**.



The reverse side shows how to create tibbles with **tibble** and to layout tidy data with **tidyr**.

OTHER TYPES OF DATA

Try one of the following packages to import other types of files

- **haven** - SPSS, Stata, and SAS files
- **readxl** - excel files (.xls and .xlsx)
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)

Save Data

Save **x**, an R object, to **path**, a file path, as:

Comma delimited file

```
write_csv(x, path, na = "NA", append = FALSE,  
          col_names = !append)
```

File with arbitrary delimiter

```
write_delim(x, path, delim = " ", na = "NA",  
            append = FALSE, col_names = !append)
```

CSV for excel

```
write_excel_csv(x, path, na = "NA", append =  
                FALSE, col_names = !append)
```

String to file

```
write_file(x, path, append = FALSE)
```

String vector to file, one element per line

```
write_lines(x, path, na = "NA", append = FALSE)
```

Object to RDS file

```
write_rds(x, path, compress = c("none", "gz",  
                                "bz2", "xz"), ...)
```

Tab delimited files

```
write_tsv(x, path, na = "NA", append = FALSE,  
          col_names = !append)
```



Read Tabular Data

- These functions share the common arguments:

```
read_*(file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"),  
      quoted_na = TRUE, comment = "", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max = min(1000,  
      n_max), progress = interactive())
```

a,b,c 1,2,3 4,5,NA	A B C 1 2 3 4 5 NA
--------------------------	--------------------------------------

a;b;c 1;2;3 4;5;NA	A B C 1 2 3 4 5 NA
--------------------------	--------------------------------------

a b c 1 2 3 4 5 NA	A B C 1 2 3 4 5 NA
--------------------------	--------------------------------------

a b c 1 2 3 4 5 NA	A B C 1 2 3 4 5 NA
--------------------------	--------------------------------------

Comma Delimited Files

```
read_csv("file.csv")
```

To make file.csv run:
write_file(x = "a,b,c\n1,2,3\n4,5,NA", path = "file.csv")

Semi-colon Delimited Files

```
read_csv2("file2.csv")
```

write_file(x = "a;b;c\n1;2;3\n4;5;NA", path = "file2.csv")

Files with Any Delimiter

```
read_delim("file.txt", delim = "|")
```

write_file(x = "a|b|c\n1|2|3\n4|5|NA", path = "file.txt")

Fixed Width Files

```
read_fwf("file.fwf", col_positions = c(1, 3, 5))
```

write_file(x = "a b c\n1 2 3\n4 5 NA", path = "file.fwf")

Tab Delimited Files

read_tsv("file.tsv") Also **read_table()**.

```
write_file(x = "a\tb\tc\n1\t2\t3\n4\t5\tNA", path = "file.tsv")
```

USEFUL ARGUMENTS

a,b,c 1,2,3 4,5,NA

Example file

```
write_file("a,b,c\n1,2,3\n4,5,NA","file.csv")  
f<- "file.csv"
```

A B C 1 2 3 4 5 NA

No header

```
read_csv(f, col_names = FALSE)
```

x y z A B C 1 2 3 4 5 NA

Provide header

```
read_csv(f, col_names = c("x", "y", "z"))
```

1 2 3 4 5 NA

Skip lines

```
read_csv(f, skip = 1)
```

A B C 1 2 3 4 5 NA

Read in a subset

```
read_csv(f, n_max = 1)
```

A B C NA 2 3 4 5 NA

Missing Values

```
read_csv(f, na = c("1", "?"))
```

Read Non-Tabular Data

Read a file into a single string

```
read_file(locale = default_locale())
```

Read each line into its own string

```
read_lines(file, skip = 0, n_max = -1L, na = character(),  
          locale = default_locale(), progress = interactive())
```

Read Apache style log files

```
read_log(file, col_names = FALSE, col_types = NULL, skip = 0, n_max = -1, progress = interactive())
```

Read a file into a raw vector

```
read_file_raw(file)
```

Read each line into a raw vector

```
read_lines_raw(file, skip = 0, n_max = -1L,  
               progress = interactive())
```



Data types

readr functions guess the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically).

A message shows the type of each column in the result.

```
## Parsed with column specification:  
## cols(  
##   age = col_integer(),  
##   sex = col_character(),  
##   earn = col_double()  
## )
```

age is an integer

sex is a character

earn is a double (numeric)

1. Use **problems()** to diagnose problems.

```
x <- read_csv("file.csv"); problems(x)
```

2. Use a **col_** function to guide parsing.

- **col_guess()** - the default
 - **col_character()**
 - **col_double()**, **col_euro_double()**
 - **col_datetime(format = "")** Also **col_date(format = "")**, **col_time(format = "")**
 - **col_factor(levels, ordered = FALSE)**
 - **col_integer()**
 - **col_logical()**
 - **col_number()**, **col_numeric()**
 - **col_skip()**
- x <- read_csv("file.csv", col_types = cols(
 A = col_double(),
 B = col_logical(),
 C = col_factor()))**

3. Else, read in as character vectors then parse with a **parse_** function.

- **parse_guess()**
 - **parse_character()**
 - **parse_datetime()** Also **parse_date()** and **parse_time()**
 - **parse_double()**
 - **parse_factor()**
 - **parse_integer()**
 - **parse_logical()**
 - **parse_number()**
- x\$A <- parse_number(x\$A)**

Tidyverse cheatsheets!

Data Transformation with dplyr :: CHEAT SHEET

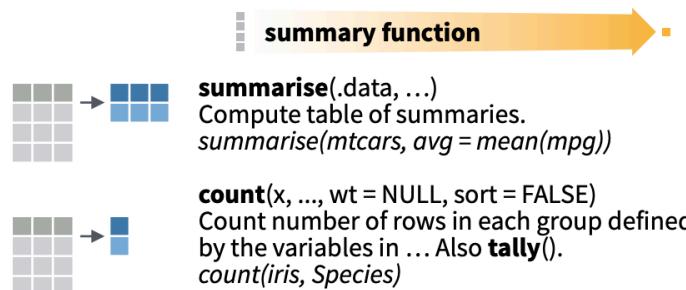


dplyr functions work with pipes and expect **tidy data**. In tidy data:



Summarise Cases

These apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).



VARIATIONS

`summarise_all()` - Apply funs to every column.
`summarise_at()` - Apply funs to specific columns.
`summarise_if()` - Apply funs to all cols of one type.

Group Cases

Use `group_by()` to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.



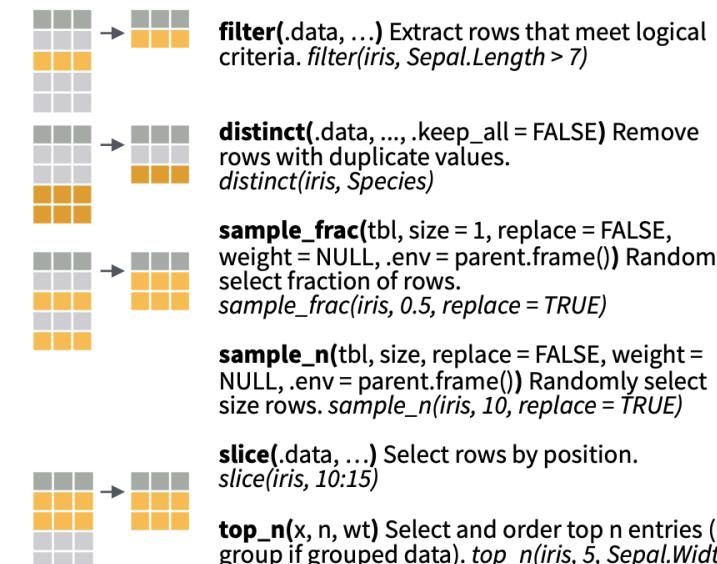
`group_by(.data, ..., add = FALSE)`
Returns copy of table grouped by ...
`g_iris <- group_by(iris, Species)`

`ungroup(x, ...)`
Returns ungrouped copy of table.
`ungroup(g_iris)`

Manipulate Cases

EXTRACT CASES

Row functions return a subset of rows as a new table.



Logical and boolean operators to use with filter()

< <= is.na() %in% | xor()
> >= !is.na() ! &

See `?base::Logic` and `?Comparison` for help.

ARRANGE CASES

An orange arrow points from a grid to a blue square icon, which then points to a blue square icon. Below it, examples are given: 'arrange(.data, ...)' (Order rows by values of a column or columns (low to high), use with `desc()` to order from high to low, e.g., `arrange(mtcars, mpg)` or `arrange(mtcars, desc(mpg))`).

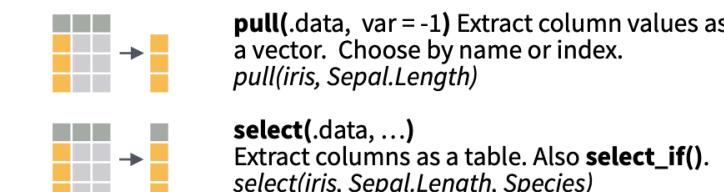
ADD CASES

An orange arrow points from a grid to a blue square icon, which then points to a blue square icon. Below it, examples are given: 'add_row(.data, ..., .before = NULL, .after = NULL)' (Add one or more rows to a table, e.g., `add_row(faithful, eruptions = 1, waiting = 1)`).

Manipulate Variables

EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.



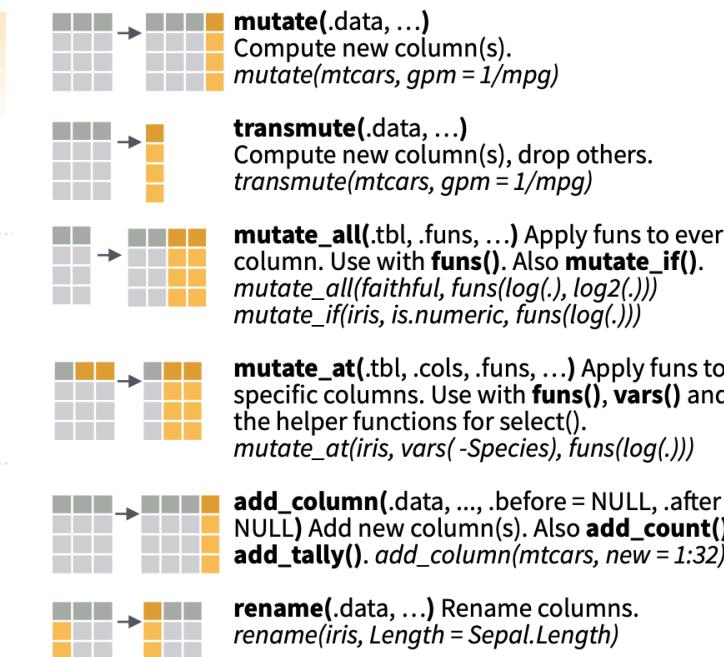
Use these helpers with `select()`, e.g., `select(iris, starts_with("Sepal"))`

`contains(match)` `num_range(prefix, range)` :, e.g., `mpg:cyl`
`ends_with(match)` `one_of(...)` -, e.g., `-Species`
`matches(match)` `starts_with(match)`

MAKE NEW VARIABLES

These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).

vectorized function →



Homework 1

Homework 1 - Data import (no Github yet)

- Open my RStudio project, which will copy it to your workspace
 - Make all changes in the script file on R Studio Cloud
 - Jake and I will be able to see your progress
 - Complete by next class
 - Questions are increasingly difficult and are meant to challenge you
 - What counts as finished? At least 4 of the questions