

ANT 6973: DATA VISUALIZATION AND EXPLORATION

# IMPORTING AND ORGANIZING DATA

# TODAY'S TOPICS

- Pipes
- File paths (here)
- Reading from flat files (`readr`)
- Reading from spreadsheets (`readxl`)

# PIPES



# PIPES



- Pass thing on left into function on the right (as the first argument)

`x %>% f()` becomes `f(x)`

`x %>% f(y)` becomes `f(x, y)`

# PIPES



- Why use this? Suppose you have a *sequence* of operations to perform on your data.

```
my_data <- 1:10
```

- *"Take the data, first exponentiate each number, then take the square root of each number, then take the mean."*

# PIPES



- Option 1: execute each function one by one, using the previous result in each subsequent function

```
result_1 <- exp(my_data)
result_2 <- sqrt(result_1)
result_final <- mean(result_2)
```

# PIPES



- Option 2: nest the functions; functions will be evaluated from inner-most to outer-most

```
result_final <- mean(sqrt(exp(my_data)))
```

# PIPES



- Option 3: chain the functions together in their natural sequence using pipes

```
result_final <- my_data %>%  
  exp() %>%  
  sqrt() %>%  
  mean()
```

- More human-readable:

*"Take the data, **then** exponentiate each number, **then** take the square root of each number, **then** take the mean."*



# PIPES



- All tidyverse packages (except for ggplot2) are designed to facilitate pipe-based workflows
- The first argument of most functions is the "data."
- ggplot2 is an older package that predates the pipe
  - It's too late to change ggplot2 now, and it has its own pipe of sorts (the "+" operator)

# PIPES



```
library("tidyverse")  
library("gapminder")  
  
count(filter(gapminder, year == 2007), continent)
```

continent	n
Africa	52
Americas	25
Asia	33
Europe	30
Oceania	2

# PIPES



```
library("tidyverse")  
library("gapminder")  
  
count(filter(gapminder, year == 2007), continent)
```

↑  
Data

continent	n
Africa	52
Americas	25
Asia	33
Europe	30
Oceania	2

# PIPES



```
library("tidyverse")
```

```
library("gapminder")
```

```
count(filter(gapminder, year == 2007), continent)
```

Step 1

continent	n
Africa	52
Americas	25
Asia	33
Europe	30
Oceania	2

# PIPES



```
library("tidyverse")  
library("gapminder")
```

```
count(filter(gapminder, year == 2007), continent)
```

← Step 2 →

continent	n
Africa	52
Americas	25
Asia	33
Europe	30
Oceania	2

# PIPES



```
library("tidyverse")  
library("gapminder")
```

```
gapminder %>%  
  filter(year == 2007) %>%  
  count(continent)
```

← Data  
← Step 1  
← Step 2

continent	n
Africa	52
Americas	25
Asia	33
Europe	30
Oceania	2

# PIPES



- Shortcut to type %>%

**Cmd** + **Shift** + **M**

(Mac)

**Ctrl** + **Shift** + **M**

(Windows)

# DATA IMPORT





# PACKAGES FOR DATA IMPORT



readr

Core part of



```
library("tidyverse")
```



readxl

Not part of core

```
library("readxl")
```

# PACKAGES FOR DATA IMPORT



haven

Import SAS, SPSS and  
STATA data files

And many more...

# PREVIOUSLY...



- Import from a .csv file using `read_csv()`

```
my_data <- read_csv("path/to/data.csv")
```

New R object where  
data will be stored

readr  
function

Path of data file relative  
to project home

# WRITE FILES



- Use the function `write_csv()` to write a data frame to disk.

```
write_csv(my_data, file = "path/my_data.csv")
```

readr  
function

Name of R  
object to export

Name & path of file you  
want to create  
note: folder must exist

# LET'S DETOUR A BIT TO LEARN MORE ABOUT PATHS...

- Go to this week's assignments on the course website.
- Download [paths-project.zip](#) and unzip the file somewhere on your computer.
- **Open the project!**
- Follow along as we complete the activity.
- *Nothing to turn in!*

20:00

# LET'S GET BACK TO DATA IMPORT



readr

- Fast way of reading rectangular data like comma separated file (csv) and tab separated files (tsv).

# FUNCTIONS



Function	Reads
<code>read_csv()</code>	Comma separated values
<code>read_csv2()</code>	Semi-colon separated values
<code>read_delim()</code>	General delimited files (e.g., " ")
<code>read_fwf()</code>	Fixed width files
<code>read_table()</code>	Space separated
<code>read_table2()</code>	Any number of whitespace characters
<code>read_tsv()</code>	Tab delimited values

# FUNCTIONS



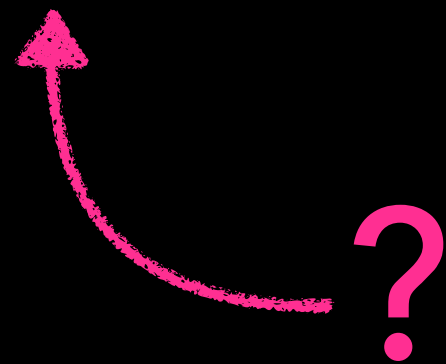
Function	Reads
<code>read_csv()</code>	Comma separated values
<code>read_csv2()</code>	Semi-colon separated values
<code>read_delim()</code>	General delimited files (e.g., " ")
<code>read_fwf()</code>	Fixed width files
<code>read_table()</code>	Space separated
<code>read_table2()</code>	Any number of whitespace characters
<code>read_tsv()</code>	Tab delimited values



# WHY USE NOT USE BASE R?



- ~10X faster than base R functions like `read.csv()`
- Returns a tibble rather than a `data.frame`



# TIBBLE

- A type of data frame used throughout the tidyverse



`read.csv()` -> `data.frame`

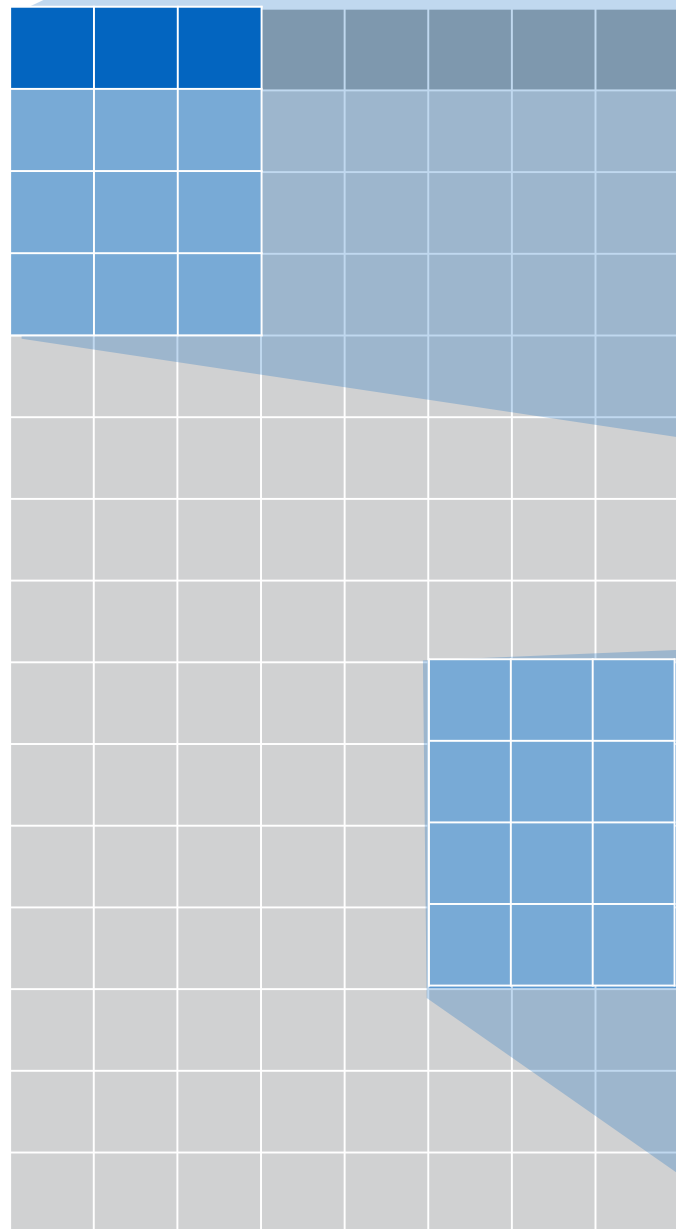
```
Console ~/Dropbox (RStudio)/RStudio/training/U-Master-the-tidyverse/0-course-developm
217 1985-10-01 -144.375 -86.5 .
218 1985-10-01 -143.125 -86.5 .
219 1985-10-01 -141.875 -86.5 .
220 1985-10-01 -140.625 -86.5 .
221 1985-10-01 -139.375 -86.5 .
222 1985-10-01 -138.125 -86.5 .
223 1985-10-01 -136.875 -86.5 .
224 1985-10-01 -135.625 -86.5 .
225 1985-10-01 -134.375 -86.5 .
226 1985-10-01 -133.125 -86.5 .
227 1985-10-01 -131.875 -86.5 .
228 1985-10-01 -130.625 -86.5 .
229 1985-10-01 -129.375 -86.5 .
230 1985-10-01 -128.125 -86.5 .
231 1985-10-01 -126.875 -86.5 .
232 1985-10-01 -125.625 -86.5 .
233 1985-10-01 -124.375 -86.5 .
234 1985-10-01 -123.125 -86.5 .
235 1985-10-01 -121.875 -86.5 .
236 1985-10-01 -120.625 -86.5 .
237 1985-10-01 -119.375 -86.5 .
238 1985-10-01 -118.125 -86.5 .
239 1985-10-01 -116.875 -86.5 .
240 1985-10-01 -115.625 -86.5 .
241 1985-10-01 -114.375 -86.5 .
242 1985-10-01 -113.125 -86.5 .
243 1985-10-01 -111.875 -86.5 .
244 1985-10-01 -110.625 -86.5 .
245 1985-10-01 -109.375 -86.5 .
246 1985-10-01 -108.125 -86.5 .
247 1985-10-01 -106.875 -86.5 .
248 1985-10-01 -105.625 -86.5 .
249 1985-10-01 -104.375 -86.5 .
250 1985-10-01 -103.125 -86.5 .
[ reached getOption("max.print") -- omitted 24974 rows ]
> |
```

`read.csv()` -> `data.frame`

```
Console ~/Dropbox (RStudio)/RStudio/training/U-Master-the-tidyverse/0-course-developm
217 1985-10-01 -144.375 -86.5 .
218 1985-10-01 -143.125 -86.5 .
219 1985-10-01 -141.875 -86.5 .
220 1985-10-01 -140.625 -86.5 .
221 1985-10-01 -139.375 -86.5 .
222 1985-10-01 -138.125 -86.5 .
223 1985-10-01 -136.875 -86.5 .
224 1985-10-01 -135.625 -86.5 .
225 1985-10-01 -134.375 -86.5 .
226 1985-10-01 -133.125 -86.5 .
227 1985-10-01 -131.875 -86.5 .
228 1985-10-01 -130.625 -86.5 .
229 1985-10-01 -129.375 -86.5 .
230 1985-10-01 -128.125 -86.5 .
231 1985-10-01 -126.875 -86.5 .
232 1985-10-01 -125.625 -86.5 .
233 1985-10-01 -124.375 -86.5 .
234 1985-10-01 -123.125 -86.5 .
235 1985-10-01 -121.875 -86.5 .
236 1985-10-01 -120.625 -86.5 .
237 1985-10-01 -119.375 -86.5 .
238 1985-10-01 -118.125 -86.5 .
239 1985-10-01 -116.875 -86.5 .
240 1985-10-01 -115.625 -86.5 .
241 1985-10-01 -114.375 -86.5 .
242 1985-10-01 -113.125 -86.5 .
243 1985-10-01 -111.875 -86.5 .
244 1985-10-01 -110.625 -86.5 .
245 1985-10-01 -109.375 -86.5 .
246 1985-10-01 -108.125 -86.5 .
247 1985-10-01 -106.875 -86.5 .
248 1985-10-01 -105.625 -86.5 .
249 1985-10-01 -104.375 -86.5 .
250 1985-10-01 -103.125 -86.5 .
[ reached getOption("max.print") -- omitted 24974 rows ]
> |
```

`read_csv()` -> `tibble`

```
Console ~/Dropbox (RStudio)/RStudio/training/U-Master-the-tidyverse/0-course-developm
> nimbus
# A tibble: 25,224 x 4
   date longitude latitude ozone
  <dtm>      <dbl>    <dbl> <chr>
1 1985-10-01 -179.375  -87.5 .
2 1985-10-01 -178.125  -87.5 .
3 1985-10-01 -176.875  -87.5 .
4 1985-10-01 -175.625  -87.5 .
5 1985-10-01 -174.375  -87.5 .
6 1985-10-01 -173.125  -87.5 .
7 1985-10-01 -171.875  -87.5 .
8 1985-10-01 -170.625  -87.5 .
9 1985-10-01 -169.375  -87.5 .
10 1985-10-01 -168.125 -87.5 .
# ... with 25,214 more rows
> |
```



A large table to display

```
# A tibble: 234 × 6  
  manufacturer model displ  
      <chr>      <chr> <dbl>  
1      audi      a4      1.8  
2      audi      a4      1.8  
3      audi      a4      2.0  
4      audi      a4      2.0  
5      audi      a4      2.8  
6      audi      a4      2.8  
7      audi      a4      3.1  
8      audi a4 quattro  1.8  
9      audi a4 quattro  1.8  
10     audi a4 quattro  2.0  
# ... with 224 more rows, and 3  
# more variables: year <int>,  
# cyl <int>, trans <chr>
```

tibble display

```
156 1999 6 auto(l4)  
157 1999 6 auto(l4)  
158 2008 6 auto(l4)  
159 2008 8 auto(s4)  
160 1999 4 manual(m5)  
161 1999 4 auto(l4)  
162 2008 4 manual(m5)  
163 2008 4 manual(m5)  
164 2008 4 auto(l4)  
165 2008 4 auto(l4)  
166 1999 4 auto(l4)  
[ reached getOption("max.print") --  
omitted 68 rows ]
```

data frame display

Dimensions

Column types

Concise view of data that fits on screen.

All other columns listed

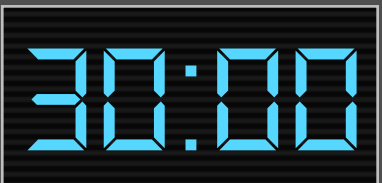
Large number of rows printed before giving up.

Many rows and variables hidden.

# ACTIVITY: OZONE



- Go to this week's assignments on the course website.
- Download `ozone.Rmd` and follow the instructions to complete the assignment.
- Submit the html report.





# AND FINALLY...



## readxl



**Jenny Bryan**  
@JennyBryan



We wanted the readxl logo to capture that special feeling you get when you hear: "We keep all our data in Excel".

[@hadleywickham](#)

♡ 532 5:54 PM - Mar 14, 2017

💬 203 people are talking about this



# data organization

organizing data in spreadsheets

My collaborators sometimes ask me, “In what form would you like the data?” My response is always, “In its current form!” If the data need to be reformatted, it’s much better for me to write a script than for them to do a bunch of cut-and-paste. I’m a strong proponent of data analysts being able to handle any data files they might receive.

But in many cases, I have to spend **a lot** of time writing scripts to rearrange the layout of the data. And how would you like your data analysts to spend their time? Reorganizing data, or really analyzing data?

Most of my collaborators enter and store their data in spreadsheets, and mostly Microsoft Excel. Before starting to enter data into a spreadsheet, it’s good to spend some time thinking about the layout. The way that you organize the data in spreadsheets can have a big impact on your data analyst’s quality of life.

This is a tutorial on that topic: *how to organize data in spreadsheets*. For complex, high-dimensional data, it may be better to use a formal database. But for many projects, spreadsheets are perfectly fine. But data in spreadsheets can be pretty and easy to work with, or they can be a sloppy mess requiring serious downstream reorganization efforts. We want to avoid the latter.

I don’t think these ideas comes naturally to anyone. So if you’re not happy with the structure of your current data files, don’t despair! And also don’t apply tedious and potentially error-prone hand-editing to revise the arrangement. Rather, apply these principles when designing the layout for your next dataset, to help make analyses easier.

- [Be consistent.](#)
- [Write dates as YYYY-MM-DD.](#)
- [Fill in all of the cells.](#)
- [Put just one thing in a cell.](#)
- [Make it a rectangle.](#)
- [Create a data dictionary.](#)
- [No calculations in the raw data files.](#)
- [Don’t use font color or highlighting as data.](#)
- [Choose good names for things.](#)
- [Make backups.](#)
- [Use data validation to avoid data entry mistakes.](#)
- [Save the data in plain text files.](#)
- [Other things to avoid.](#)
- [Other resources](#)



# UGH...



dirty\_data

Search Sheet

Tell me what you want to do

Share

Home Insert Page Layout Formulas Data Review View

Paste Cut Copy Format

Calibri (Body) 11 A A

B I U

Wrap Text

Merge & Center

General

\$ %

Conditional Formatting

Format as Table

Normal Bad Good Neutral

Insert Delete Format

AutoSum Fill Clear

Sort & Filter

D13

	A	B	C	D	E	F	G	H	I	J	K
1	First Name	Last Name	Employee Status	Subject	Hire Date	% Allocated	Full time?	do not edit! --->	Certification	Certification	Certification
2	Jason	Bourne	Teacher	PE	39690	75%	Yes		Physical ed	Theater	
3	Jason	Bourne	Teacher	Drafting	39690	25%	Yes		Physical ed	Theater	
4	Alicia	Keys	Teacher	Music	37118	100%	Yes		Instr. music	Vocal	
5	Ada	Lovelace	Teacher	#REF!	27515	100%	Yes		PENDING	Comp	
6	Desus	Nice	Administration	Dean	41431	100%	Yes		PENDING		
7	Chien-Shiung	Wu	Teacher	Physics	11037	50%	Yes		Science 6-12	Physics	
8	Chien-Shiung	Wu	Teacher	Chemistry	11037	50%	Yes		Science 6-12	Physics	
9											
10	James	Joyce	Teacher	English	32994	50%	No			English 6-12	
11	Hedy	Lamarr	Teacher	Science	27919	50%	No		PENDING		
12	Carlos	Boozer	Coach	Basketball	42221	#N/A	No		Physical ed		
13	Young	Boozer	Coach		34700	#N/A	No			Political sci.	
14	Micheal	Larsen	Teacher	English	40071	80%	No		Vocal music	English	
15											
16											
17											
18											
19											
20											

Sam: Put alternate certifications in columns to right

Sheet1

Ready

200%

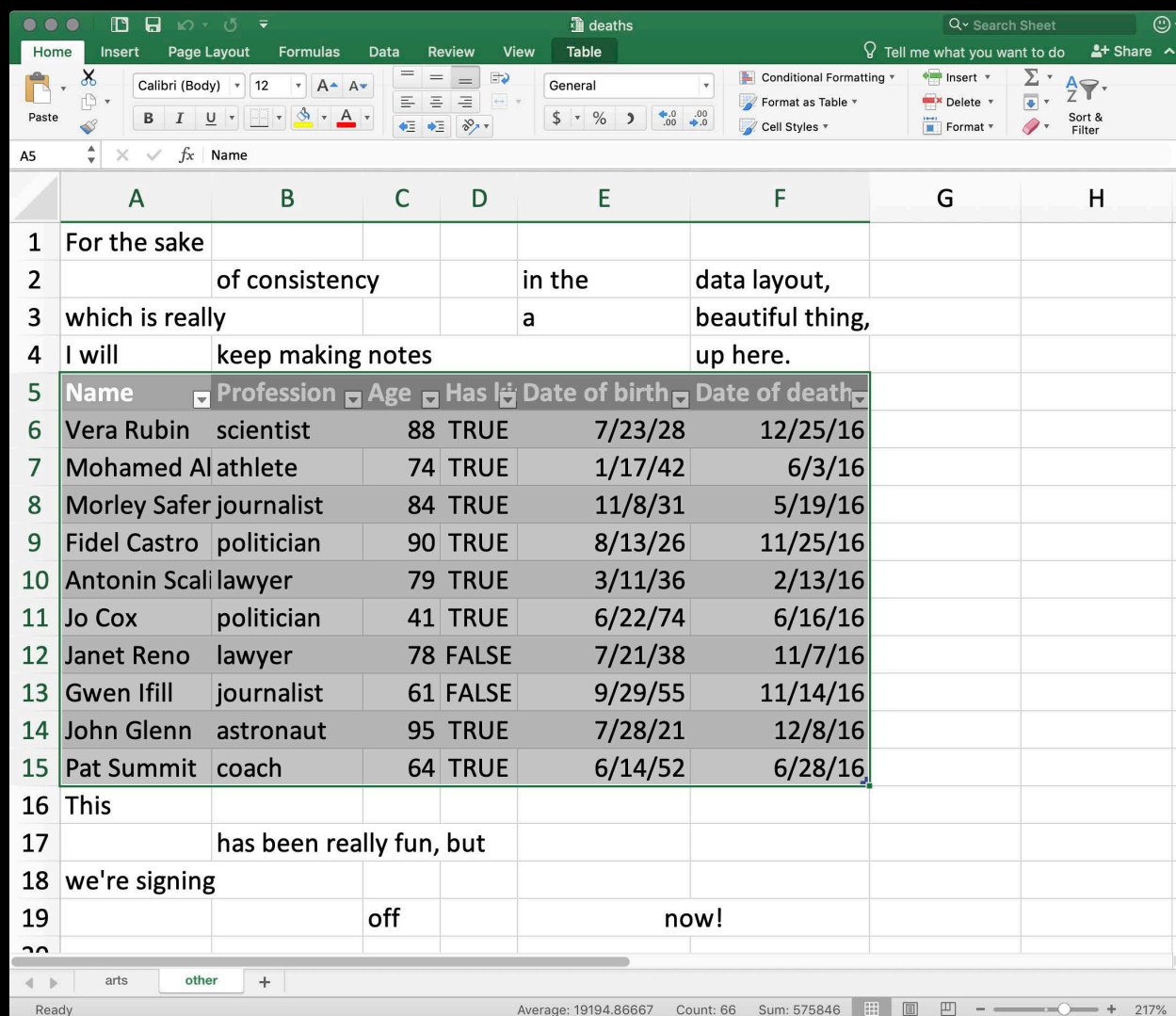
# READXL: A PRIMER



	A	B	C	D	E	F	G	H
1	For the sake							
2		of consistency			in the	data layout,		
3	which is really				a	beautiful thing,		
4	I will	keep making notes				up here.		
5	Name	Profession	Age	Has l	Date of birth	Date of death		
6	Vera Rubin	scientist	88	TRUE	7/23/28	12/25/16		
7	Mohamed Al	athlete	74	TRUE	1/17/42	6/3/16		
8	M							
9	Fi							
10	Ar							
11	Jo Cox	politician	41	TRUE	6/22/74	6/16/16		
12	Janet Reno	lawyer	78	FALSE	7/21/38	11/7/16		
13	Gwen Ifill	journalist	61	FALSE	9/29/55	11/14/16		
14	John Glenn	astronaut	95	TRUE	7/28/21	12/8/16		
15	Pat Summit	coach	64	TRUE	6/14/52	6/28/16		
16	This							
17		has been really fun, but						
18	we're signing							
19		off			now!			



readxl is all about the data rectangle



	A	B	C	D	E	F	G	H
1	For the sake							
2		of consistency		in the		data layout,		
3	which is really			a		beautiful thing,		
4	I will	keep making notes				up here.		
5	Name	Profession	Age	Has	Date of birth	Date of death		
6	Vera Rubin	scientist	88	TRUE	7/23/28	12/25/16		
7	Mohamed Al	athlete	74	TRUE	1/17/42	6/3/16		
8	Morley Safer	journalist	84	TRUE	11/8/31	5/19/16		
9	Fidel Castro	politician	90	TRUE	8/13/26	11/25/16		
10	Antonin Scali	lawyer	79	TRUE	3/11/36	2/13/16		
11	Jo Cox	politician	41	TRUE	6/22/74	6/16/16		
12	Janet Reno	lawyer	78	FALSE	7/21/38	11/7/16		
13	Gwen Ifill	journalist	61	FALSE	9/29/55	11/14/16		
14	John Glenn	astronaut	95	TRUE	7/28/21	12/8/16		
15	Pat Summit	coach	64	TRUE	6/14/52	6/28/16		
16	This							
17		has been really fun, but						
18	we're signing							
19		off			now!			



```
read_excel("deaths.xlsx", skip = 4)
read_excel("deaths.xlsx", n_max = 15)
read_excel("deaths.xlsx", skip = 4, n_max = 10)
read_excel("deaths.xlsx", range = "A5:F15")
read_excel("deaths.xlsx", range = cell_rows(5:15))
read_excel("deaths.xlsx", range = cell_cols("A:F"))
```



RStudio has an Excel helper than works with readxl

The screenshot displays the RStudio environment. The 'File' menu is open, and the 'Import Dataset' option is selected, leading to a submenu where 'From Excel...' is chosen. The console window shows the following R code and output:

```
# ... with 132 more rows
> gapminder %>%
+   filter(year == 2007) %>% count(continent)
# A tibble: 5 x 2
  continent      n
  <fct>         <int>
1 Africa         52
2 Americas        25
3 Asia            33
4 Europe          30
5 Oceania          2
>
```

The bottom right pane shows a file explorer with the following files:

File Name	Size	Modified
spreadsheets		
spreadsheets-solutions.html	730 KB	Mar 6, 2019, 3:53 PM
spreadsheets-solutions.Rmd	5.4 KB	Mar 6, 2019, 4:15 PM
spreadsheets.html	728.7 KB	Mar 6, 2019, 3:46 PM
spreadsheets.Rmd	3.4 KB	Mar 6, 2019, 3:56 PM



Environment

History

Connections

Console

Terminal

R Markdown

Global Env

Data

gm\_2007

Values

my\_data

### Import Excel Data

File/Url:

~/OneDrive - University of Texas at San Antonio/Teaching/Data Visualization/activities/ant6973-activities/08\_data-import/spreadsheets/deaths.xlsx

Browse...

Data Preview:

Name (character)	Profession (character)	Age (numeric)	Has kids (double)	Date of birth (date)	Date of death (date)
David Bowie	musician	69	1	Character	2016-01-10
Carrie Fisher	actor	60	1	Numeric	2016-12-27
Chuck Berry	musician	90	1	Date	2017-03-18
Bill Paxton	actor	61	1	Include	2017-02-25
Prince	musician	57	1	Skip	2016-04-21
Alan Rickman	actor	69	0		2016-01-14
Florence Henderson	actor	82	1	1934-02-14	2016-11-24
Harper Lee	author	89	0	1926-04-28	2016-02-19
Zsa Zsa Gabor	actor	99	1	1917-02-06	2016-12-18
George Michael	musician	53	0	1963-06-25	2016-12-25

Previewing first 50 entries.

Import Options:

Name:	deaths	Max Rows:	10	<input checked="" type="checkbox"/> First Row as Names
Sheet:	Default	Skip:	4	<input checked="" type="checkbox"/> Open Data Viewer
Range:	A1:D10	NA:		

Code Preview:

```
library(readxl)
deaths <- read_excel("08_data-import/spreadsheets/deaths.xlsx",
  col_types = c("text", "text", "numeric",
    "numeric", "date", "date"), skip = 4,
  n_max = 10)
View(deaths)
```

? Reading Excel files using readxl

Import

Cancel



```
~/OneDrive - University of Texas at San Antonio/Teaching/Data Visualization/activities/ant6973-activities - RStudio
```

Go to file/function | Addins | ant6973-activities

Console | Terminal x

```
> read_excel(here("08_data-import/spreadsheets", "deaths.xlsx"),
+           sheet = "other", cell_rows(5:15))
# A tibble: 10 x 6
```

	Name	Profession	Age	`Has kids`	`Date of birth`	`Date of death`
	<chr>	<chr>	<dbl>	<lgl>	<dtm>	<dtm>
1	Vera Rubin	scientist	88	TRUE	1928-07-23 00:00:00	2016-12-25 00:00:00
2	Mohamed Ali	athlete	74	TRUE	1942-01-17 00:00:00	2016-06-03 00:00:00
3	Morley Safer	journalist	84	TRUE	1931-11-08 00:00:00	2016-05-19 00:00:00
4	Fidel Castro	politician	90	TRUE	1926-08-13 00:00:00	2016-11-25 00:00:00
5	Antonin Scalia	lawyer	79	TRUE	1936-03-11 00:00:00	2016-02-13 00:00:00
6	Jo Cox	politician	41	TRUE	1974-06-22 00:00:00	2016-06-16 00:00:00
7	Janet Reno	lawyer	78	FALSE	1938-07-21 00:00:00	2016-11-07 00:00:00
8	Gwen Ifill	journalist	61	FALSE	1955-09-29 00:00:00	2016-11-14 00:00:00
9	John Glenn	astronaut	95	TRUE	1921-07-28 00:00:00	2016-12-08 00:00:00
10	Pat Summit	coach	64	TRUE	1952-06-14 00:00:00	2016-06-28 00:00:00

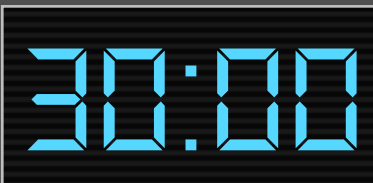
```
>
```

readxl is good at guessing column types

# ACTIVITY: SPREADSHEETS



- Go to this week's assignments on the course website.
- Download `spreadsheets.Rmd` and follow the instructions to complete the assignment.
- Submit the html report.





# ACKNOWLEDGEMENTS

- Some ideas, examples, and figures from RStudio webinars, which are licensed CC by SA.