

Introduction to R

CRDDS Summer
Research Data
Bootcamp

May 17, 2023

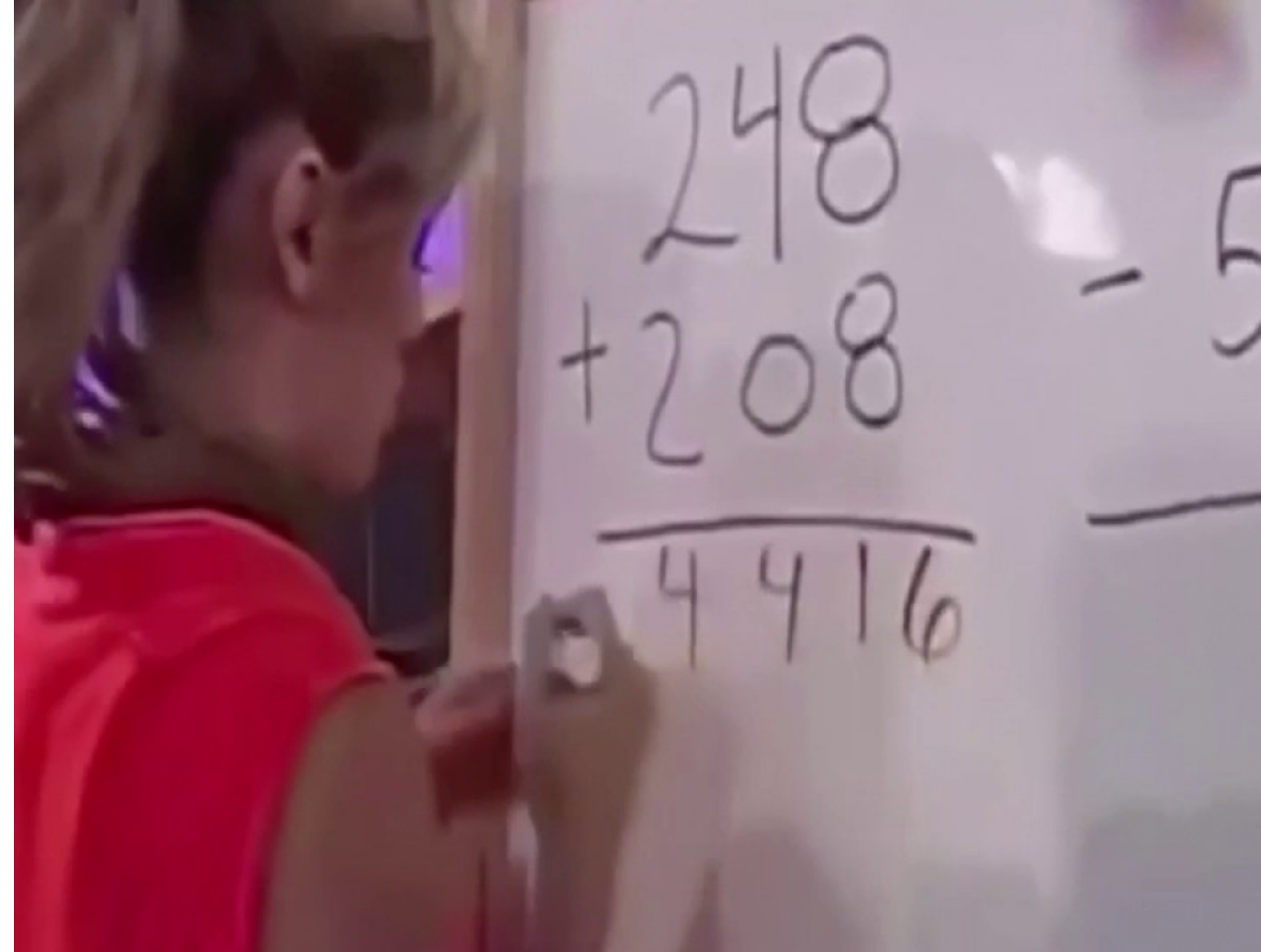
Scott Nordstrom

(Scott.Nordstrom@Colorado.edu,
scottwatsonnordstrom@gmail.com)



Why use R?

- Less error-prone
- More efficient and faster
- Handle large, disparate datasets
- Reproducible
 - Others can re-create your work
 - You can re-create your work
 - Reliable – should give same answer each time



R vs. RStudio



- R is the language, RStudio is software
- Analogy: if R was English, then RStudio would be Microsoft Word
- Further torturing the analogy: R is English, Rstudio is Microsoft Word, a .R file ("script") is a poem



RStudio: what's going on here?

The screenshot shows the RStudio environment with four main panels. The top-left panel is the Editor window, containing an R script. The bottom-left panel is the Console window, showing the execution of the script. The top-right panel displays the Workspace and History, and the bottom-right panel shows the R Documentation for the `lm` function. Annotations include a red box around the Editor window, a yellow box around the Console window, and a black box with arrows pointing to the Workspace and Documentation panels.

Editor window

Your script will open in here.
You can run whole scripts at once or just parts of scripts.

```
1  
2 rm(list = ls())  
3 N <- 1000  
4 u <- rnorm(N)  
5 x1 <- -2 + rnorm(N)  
6 x2 <- 1 + x1 + rnorm(N)  
7 y <- 1 + x1 + x2 + u  
8 r1 <- lm(y ~ x1 + x2)  
9  
10 |
```

Console window

You can run one or more commands here.
Output will be printed here.

```
> rm(list = ls())  
> N <- 1000  
> u <- rnorm(N)  
> x1 <- -2 + rnorm(N)  
> x2 <- 1 + x1 + rnorm(N)  
> y <- 1 + x1 + x2 + u  
> r1 <- lm(y ~ x1 + x2)  
> |
```

Workspace

Variable	Value
N	1000
r1	lm[12]
u	numeric[1000]
x1	numeric[1000]
x2	numeric[1000]
y	numeric[1000]

R Documentation: Fitting Linear Models

Description

`lm` is used to fit linear models. It can be used to carry out regression, single stratum analysis of variance and analysis of covariance (although `av` may provide a more convenient interface for these).

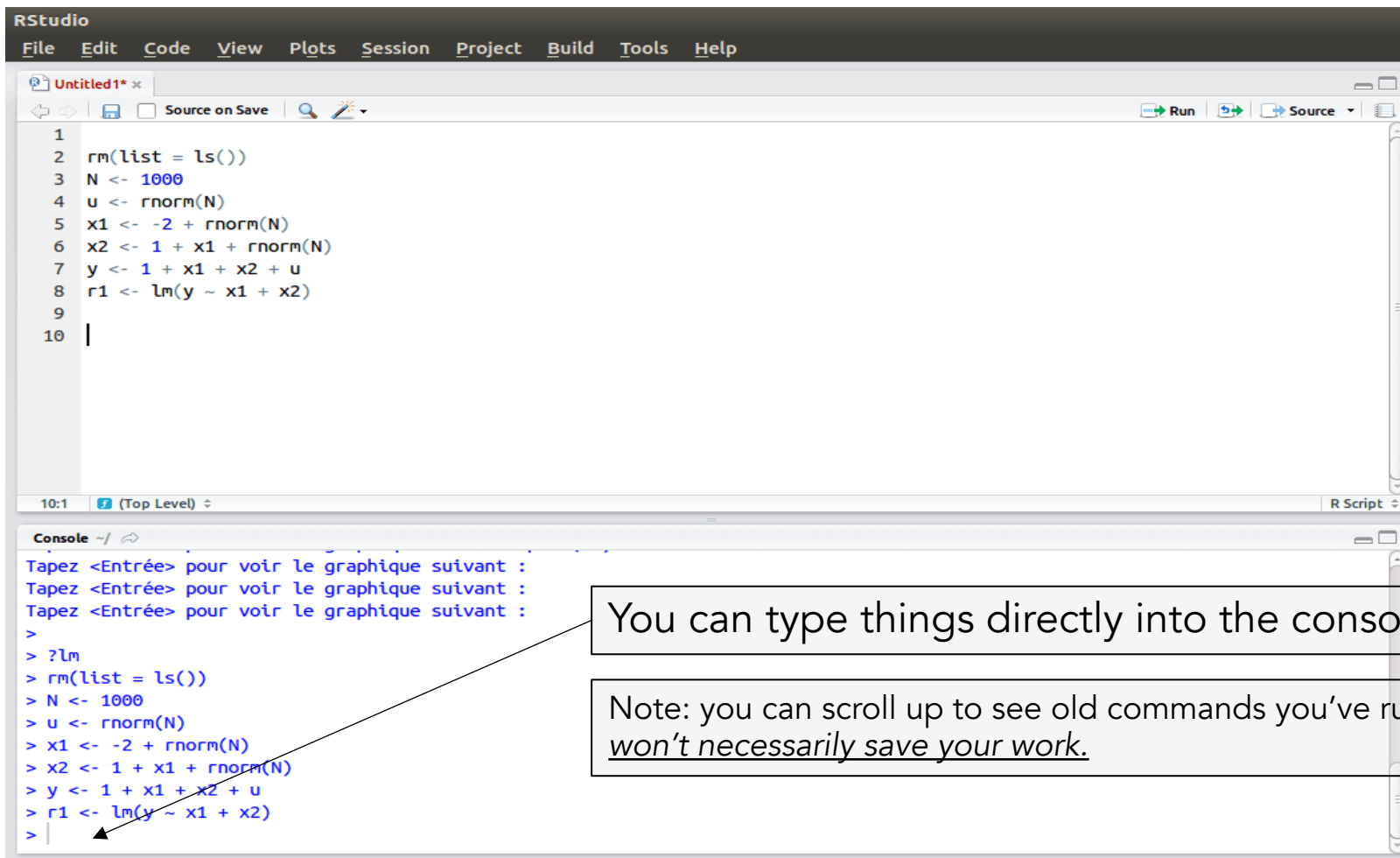
Usage

```
lm(formula, data, subset, weights,  
    method = "qr", model = TRUE, x =  
    singular.ok = TRUE, contrasts =
```

Arguments

Let's ignore these for now.

Running R commands in the console



The screenshot shows the RStudio environment. The top menu bar includes File, Edit, Code, View, Plots, Session, Project, Build, Tools, and Help. Below the menu is a toolbar with icons for running and saving. The main editor window, titled 'Untitled1*', contains the following R code:

```
1  
2 rm(list = ls())  
3 N <- 1000  
4 u <- rnorm(N)  
5 x1 <- -2 + rnorm(N)  
6 x2 <- 1 + x1 + rnorm(N)  
7 y <- 1 + x1 + x2 + u  
8 r1 <- lm(y ~ x1 + x2)  
9  
10 |
```

The console window at the bottom shows the execution of the same code, with the prompt '>' at the start of each line. The text 'Tapez <Entrée> pour voir le graphique suivant :' is repeated three times in blue. The console output is as follows:

```
>  
> ?lm  
> rm(list = ls())  
> N <- 1000  
> u <- rnorm(N)  
> x1 <- -2 + rnorm(N)  
> x2 <- 1 + x1 + rnorm(N)  
> y <- 1 + x1 + x2 + u  
> r1 <- lm(y ~ x1 + x2)  
> |
```

An arrow points from the text box below to the prompt '>' in the console.

You can type things directly into the console and hit "Enter"

Note: you can scroll up to see old commands you've run, but otherwise, this won't necessarily save your work.



C'mon,
do something...



Arithmetic and logic in R

- R (and any programming language) can do simple arithmetic
- Try running the following commands in the console:

`2 + 2`
`[1] 4`

I'll use this font throughout the workshop denote R code

The font is "Monaco" if you ever want to use it – it's the same as the default in RStudio!

This is "returned" by the command – running commands usually prints the output in the console.

`(2 + 2)^2`

Keep your order of operations (PEMDAS) in mind!

`19^2 < 4^4`

This is what's called a "logical statement" or a "boolean statement" – like asking a question where the answer is TRUE or FALSE

Running R commands from the editor

RStudio

File Edit Code View Plots Session Project Build Tools Help

or "ctrl + enter/return"

Run

1
2 rm(list = ls())
3 N <- 1000
4 u <- rnorm(N)
5 x1 <- -2 + rnorm(N)
6 x2 <- 1 + x1 + rnorm(N)
7 y <- 1 + x1 + x2 + u
8 r1 <- lm(y ~ x1 + x2)
9
10

You can also type commands into a script, then highlight them and click "Run" or highlight them and hit "ctrl + Enter"

"Running" code will send it to the console.

This will save your commands to re-run later, but remember, what you type into a script won't run less you run it!

10:1 (Top Level)

Console

Tapez <Entrée> pour voir le graphique suivant :
Tapez <Entrée> pour voir le graphique suivant :
Tapez <Entrée> pour voir le graphique suivant :
>
> ?lm
> rm(list = ls())
> N <- 1000
> u <- rnorm(N)
> x1 <- -2 + rnorm(N)
> x2 <- 1 + x1 + rnorm(N)
> y <- 1 + x1 + x2 + u
> r1 <- lm(y ~ x1 + x2)
>

You can type things directly into the console and hit "Enter"

Note: you can scroll up to see old commands you've run, but otherwise, this won't necessarily save your work.

How I (Scott) do it:

- Use the console for trying stuff out
- Do what I want to save in the editor



C'mon,
do something...



"Hanging" (incomplete) commands

Common source of error is entering incomplete commands.

When the console gets an incomplete command, *it waits for you to close the command* before running. The "+" is like an invitation to finish the command.

You can fix this by finishing the command.

This happens commonly with arithmetic (e.g., +, /), parentheses and brackets, and quotes.

The diagram illustrates a console session with two parts. In the first part, a prompt ">" is followed by "2 +". A blue arrow points from the text "Your command, in the editor" to the "2 +". Below this, a new line starts with a "+" followed by a vertical bar "|". A blue arrow points from the text "No >, line starts with +, blinking cursor, command did not run." to the "+" and the bar. In the second part, the prompt ">" is followed by "2 +". Below this, a new line starts with "+" followed by "2". A blue arrow points from the text "Put in the rest of your command here." to the "2". Below that, the prompt "[1]" is followed by "4". At the bottom, the prompt ">" is followed by a vertical bar "|". A blue arrow points from the text "Now you're good to enter another command!" to the bar.

Your command, in the editor

No >, line starts with +, blinking cursor, command did not run.

Put in the rest of your command here.

Now you're good to enter another command!

Variables: saving objects in memory

`x <- 2 + 2`

`x`

`(x <- 2 + 2)`

This will save the output (right of the arrow) into memory under the name `x`

Running just the variable name itself will print its value

(useful to make sure the output is what you expected it to be!)

Little known trick – wrapping this statement in parentheses will print `x` to the console too!

Useful rules for naming variables:

- Can't include a space
- Can't start with numbers
- Can include underscore (`_`) and period (`.`)
- Names are case sensitive

Widely followed conventions:

- Use informative names!
- Start with lowercase letters
- camelCase or underscores for names with multiple words
 - E.g., `myData` or `my_data` – easier to read than `mydata`

<- VS. =

- Both are “assignment” operators – they do the same thing
 - <- has extra functionality in some high-level circumstances
- Most people use <- but I prefer to use =
 - = is one keystroke while <- is three!



$$x = 2 + 2$$

$$y = 3 * 18$$

$$z = y / x$$

You can plug variables into commands!

Variables: when they stay and when they go

Declare variables and they stay in memory unless

- You overwrite them
- You remove them (`rm()`)
- You restart RStudio (or the session)

```
x = -20  
x = 4*4  
x  
  
rm(x)  
x
```

See your variables with the command `ls()`

If you're running a script and your R session crashes, all your variables will be wiped out...

- (but if you save your script, you can just run everything again!)

Comment with # – hard

If you include the # (pound sign) in
R will ignore (not run) everything after

`2 + 2 + 2 # + 2` *What will this return*

`2 + 2 + # 2 + 2` *What will happen with*

```
101 ### Define a global variable for max age to use in model:
102 max.age = 6 # for now set to 6 (years 2013 - 2019)
103
104 ### Define variables for model
105 vars = c(paste0('ld.age', 0:max.age), "lfCt.age6")
106
107 length(vars)
108 vars
109
110 ### Defining the graphical structure (predecessor)
111 # (I think this should just be linear?)
112 pred = 0:(max.age+1)
113
114 # Check to make sure this works
115 rbind(succ = vars, pred = c("initial", vars)[pred + 1])
116
117 ### Defining families
118 # Every living during variable should be bernoulli
119 # Assume that leaf count is zero-truncated Poisson (check with group)
120 fam = c(rep(1, (length(vars)) - 1), 3)
121
122 ### Reshaping data frame:
123
124 # It's already in long form
125 head(t2)
```

These are great to include to “narrate” scripts
useful to someone else looking at your code,
but also useful to future you looking at your code!

Functions

Commands that have parentheses after them are functions.

Most functions have an **input(s)** and **output**.

The inputs are often called arguments.

(you can write your own functions!)

`paste(..., sep = ' ')`

↑
Characters you want to paste together to here.

This function will combine a vector of characters into one object, with each object in the vector separated by whatever you give as the argument "`sep`"

The "argument" goes here – it's what you want to take the cosine of

`cos()`

`cos(0)`

We usually say the function "returns" some output – in this case it would return "1".

Getting help

? Followed by the name of a function will pull up the help/documentation page for that function

? `sqrt`

Other data “types”

So far we have dealt only with integers.

Unsurprisingly, R can also handle non-integer numbers.

```
sqrt(15)  
[1] 3.872983
```

Other types:

- “character” – wrapped in quotes (single or double)

```
“five”
```

```
five <- “five”  
type(five)
```

- “logical” – TRUE or FALSE

```
5 > 4
```

```
[1] TRUE
```

```
five.four <- 5 > 4
```


Store multiple objects in **vectors**

Make a vector with the command `c()`

```
brady_sb_wins <- c(2001, 2003, 2004, 2014, 2016, 2018, 2020)
```

```
brady_sb_wins <- c('xxxvi', 'xxxviii', 'xxxix', 'xlix', 'li',  
'liii', 'lv')
```

Vector elements must all be of the same type.



```
ex_vec <- c(5, 6, 'red')
```

```
ex_vec
```

```
[1] "5"      "6"      "red"
```

*all elements were
coerced into characters*

“Indexing” - picking out a slice or piece

What if I want to pick out just one (or a couple) observations?

For a vector: use brackets []

letters

```
[1] "a" "b" "c" "d" "e" "f" "g" "h" "i" "j" "k" "l" "m" "n" "o" "p" "q" "r" "s" "t" "u" "v" "w"  
"x" "y" "z"
```

letters[5]

```
[1] "e"
```

The [5] here says “pick out just the fifth element”

Indexing a vector

Say now I want to pick out multiple elements (but not all of them)

Make a vector of the elements you want to pick out using `c()` and stick that in brackets.

Pick out elements, 5, 22, 15, ...

```
> letters[c(5, 22, 15, 12, 21, 20, 9, 15, 14)]  
[1] "e" "v" "o" "l" "u" "t" "i" "o" "n"
```

Add to a vector:

```
usa.colors <- c('red', 'white', 'blue')
```

```
azeri.colors <- c(usa.colors, 'green')
```

```
azeri.colors
```

```
[1] "red"    "white"  "blue"   "green"
```

```
lebanese.colors <- azeri.colors[-3]
```

```
lebanese.colors
```

```
[1] "red"    "white"  "green"
```



"Minus"
indices to
remove them
(you can do
this with
multiple!, e.g.,
try `-c(3, 4)`)

Modifying vectors

```
squares <- c(1, 4, 9, 16, 25)
```

```
sqrt(squares[1])
```

```
sqrt(squares[2])
```

```
...
```

```
square.roots <- sqrt(squares)
```

```
square.roots
```

```
[1] 1 2 3 4 5
```

Many functions in R are written to be
"vectorized", i.e., can perform operations
on each element independently



Break!

Data frames

2D object meant for handling datasets (think: spreadsheets)

- Rows: typically correspond to one "data point"

E.g., the "cars" data frame:

```
> cars
  speed dist
1     4    2
2     4   10
3     7    4
4     7   22
5     8   16
```

Columns named "speed" and "dist"
hold data for one observation.

Probably variables in your analysis!

- Columns: can be numbers or characters

Rows correspond to observations (of a car?).

R has a lot of built-in features meant to work with data frames

Indexing a data frame

Data frames have rows and columns

You'll still use brackets, but now there are *two* indices

`[rows, columns]`

```
cars[1,1]
```

```
[1] 4
```

Picks out first row and first column

```
cars[c(5, 7, 10), 2]
```

```
[1] 16 18 17
```

Picks out rows 5, 7, 10 and second column

Indexing a data frame

But remember...

R has a lot of built-in features meant to work with data frames

Pick out one column of a data frame with the dollar sign (\$)

`cars$speed`

`cars$speed[5]`

Pick out multiple columns using character vectors

`cars[,c("dist", "speed")]`

*Note that indexing is order specific!
What did this do?*

More fun with indexing!

You can index with variables you have made!

```
fast_cars = c(47, 48, 49)
```

```
cars[fast_cars, ]
```

	speed	dist
47	24	92
48	24	93
49	24	120



You can index with logical statements!

What is this command doing?

```
cars[cars$speed < 15,]
```

`==` in logical statements – don't confuse with `=`

```
cars[cars$speed == 15,]
```

```
cars[cars$speed != 15,]
```

Exclamation mark *negates* logical statements

The dreaded error...

Error – some part of your code did not execute

- Error messages – these are informative! Read them!
- You can also google error messages!

“Debugging” – fixing “bugs” in your code



Note – just because you don't get an error message doesn't mean your code runs *as you want it!*



Donald J. Trump ✓
@realDonaldTrump

Thank you Kanye, very cool!

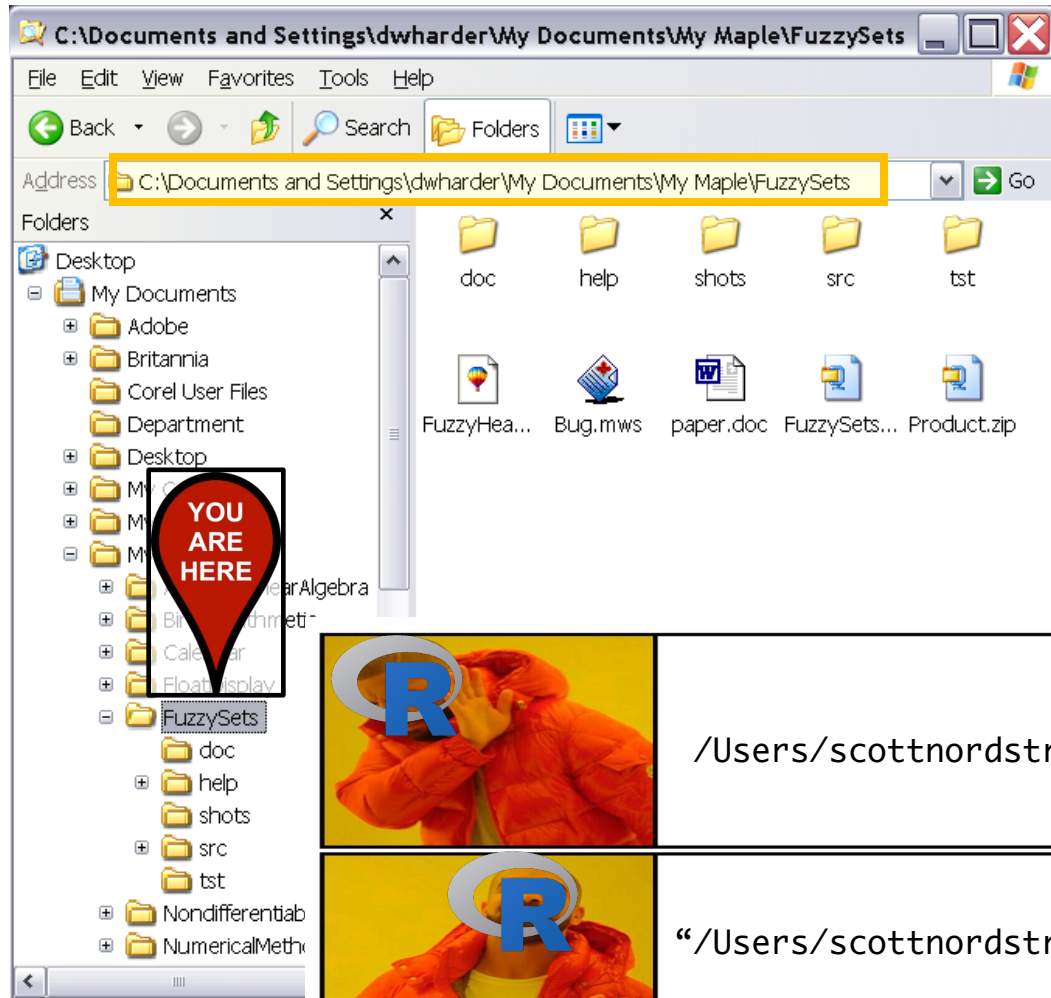
KANYE WEST ✓ @kanyewest

You don't have to agree with trump but the mob can't make me not love him. We are both dragon energy. He is my brother. I love everyone. I don't agree with everything anyone does. That's what makes us individuals. And we have the right to independent thought.

All of that is very cool.

How do I analyze my own data though...?

Directories – where am I?



R sessions are always “parked” in some directory on your computer.

Where it is parked is called your “*working directory*.”

Check your working directory: `getwd()`

Set your working directory: `setwd()`

Directories (and file names) should be in quotes!



`/Users/scottnordstrom/`



`“/Users/scottnordstrom/”`

Reading in files

`read.csv()`

CSV = “comma separated value” – R can not read in .xsl or .xlsx files
(unless you install extra packages)

`myData <- read.csv(filename)`

You can always read in the *full* path to your file

```
read.csv("/Users/scottnordstrom/Teaching/r_crdds_2023-05/wos/ClimateAndArt1.csv")
```

If your csv is in your working directory, you can just put the filename

```
setwd("/Users/scottnordstrom/Teaching/r_crdds_2023-05/wos")  
read.csv("ClimateAndArt1.csv")
```

If your csv is in a *subdirectory* of your working directory, you can put the path to the file

```
setwd("/Users/scottnordstrom/Teaching/r_crdds_2023-05")  
read.csv("wos/ClimateAndArt1.csv")
```

Possibly useful function:
`file.exists()`

Read in the first ClimateAndArt CSV

```
# (remember to set your working directory first)
climateArt1 <- read.csv('ClimateAndArt_01.csv')
```

Inspecting the data frame:

`dim(climateArt1)` will tell us how many rows and columns

`names(climateArt1)` will show names of each column

Other useful functions:

`head(climateArt1)` will print out first several rows (you can guess what `tail()` does)

`str(climateArt1)` will show us the **type** stored in each column

Exploring the data frame

What publication types are there in the data frame?

```
climateArt1$Publication.type (this will print the whole thing!)  
unique(climateArt1$Publication.type)
```

How many of each? – **table()**
`table(climateArt1$Publication.type)`

Combining data frames with `rbind()`

`rbind()` will combine data frames into one. ('r' for 'row')

Requires all data frames to have identical column names

```
climateArt <- rbind(  
  read.csv(...),  
  read.csv(...),  
  ...  
)
```

Can you guess what `cbind()` does?

Can you guess what must be true for `cbind()` to work?

There are other, sleeker ways to do this – if you are curious, ask me!

Tidy data: working with the tidyverse

Developers have made handy packages for handling data-types.

These are called the tidyverse.

We will get started with the tidyverse after the break.

For now, make sure the following lines work for you:

```
library(dplyr)
library(tidyr)
library(ggplot2)
```





Break!

On packages ("libraries")

"Base" R is a collection of functions that run on their own.

Install packages with
`install.packages()`

But sometimes, people figure out ways to do things better, faster, more neatly, etc.

`install.packages('nic')`

This is a package with "nature inspired color palettes"

They bundle this code into external "packages" that you can install and use.

You only need to install packages *once* (unless you update your version of R)

Don't ever include `install.packages()` in a script! It will just re-install the package, wasting time.

Manipulating and preparing data with the tidyverse

- Tidyverse is a collection of packages for manipulating data (and other things)
- “Base” R: confusing, inconsistent hodgepodge of functions
- Tidyverse (ideally): more consistent, coherent structure/organization



Consistent structure means a little bit of knowledge goes a long way

Do I need to use the tidyverse?

No.

Much of what the tidyverse is capable of can be done in base R.

But, it might be slower, messier, more complicated, etc.

Tidyverse pros: powerful, widely used, common syntax means learning new things is easier once you've seen enough

Tidyverse cons: learning curve, occasionally changing (old code might "break")



Getting started with tidyverse

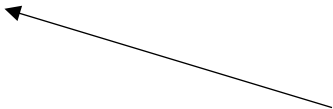
To make the functions in a package accessible, use the `library()` command to load the package.

Today, we'll use these tidyverse packages: `tidyr`, `dplyr`, `ggplot2`

```
library(dplyr)
library(tidyr)
library(ggplot2)
```

It generally *is* a good idea to include calls to `library()` in your scripts, usually at the beginning.

You can also say `library(tidyverse)`, which will load *all* of the tidyverse packages at once.



Fundamental tidyverse concept: pipe

Pipe: %>%

```
object %>% function()
```

Like saying “take *object* *and then* put it into `function()`”

More formally: the *object* on the left-hand side gets “piped” in as the first argument in `function()`

```
object %>% function(argument2, argument 3)
```

If more than one argument to `function()`, pass those in after

Simple pipe examples

```
5 %>% sqrt()
```

```
sqrt(5)
```

```
mtcars$cyl %>% character()
```

```
character(mtcars$cyl)
```

```
mtcars$cyl %>%  
  nrow() %>%  
  sqrt()
```

```
n <- nrow(mtcars$cyl)  
sqrt(n)
```

or

```
sqrt(nrow(mtcars$cyl))
```

No temporary objects

Easy to read

Makes unnecessary variables

Harder to read



Piping usefulness: stringing together operations

```
595 therm %>%  
596   filter(Plot %in% 48) %>%  
597   filter(!Tag %in% 1117) %>%  
598   select(Date, Tag, Infl_spread) %>%  
599   spread(Date, Infl_spread)  
600 # 1068 is plusible...  
601 therm %>%  
602   filter(Plot %in% 48) %>%  
603   filter(!Tag %in% 1117) %>%  
604   select(Date, Tag, Infl_done) %>%  
605   spread(Date, Infl_done)  
606 # 1068 is the only realistic option
```

Neatly, cleanly perform multiple operations on data frame

No temporary objects made

Easy to read (if you know what you're looking for)

```
640 therm %>%  
641   filter(Plot %in% 48) %>%  
642   filter(!Tag %in% 1117) %>%  
643   select(Date, Tag, Infl_spread, Infl_done) %>%  
644   unite(col = Infl, c(Infl_spread, Infl_done), sep = '_') %>%  
645   spread(Date, Infl)
```

Another analogy for piping: cake baking



Cake recipe:

- Take flour
- Add eggs, oil, water
- Mix with spoon for two minutes
- Bake at 350 degrees F for 35 minutes
- Let cool

Base R:

```
dough <- add(flour, oil, water)
batter <- mix(dough, utensil =
              'spoon', time = 2)
cake <- bake(batter, temp = 350,
             unit = 'F', time = 35)
cake <- let_cool(cake)
```

dough, batter are made once, never used again

Another analogy for piping: cake baking



Cake recipe:

- Take flour
- Add eggs, oil, water
- Mix with spoon for two minutes
- Bake at 350 degrees F for 35 minutes
- Let cool

Base R:

```
cake <- bake(mix(add(flour, oil, water),  
  utensil = 'spoon', time = 2)  
  temp = 350, unit = 'F',  
  time = 35)  
)
```

```
cake <- let_cool(cake)
```

this really hard to read!

Another analogy for piping: cake baking



Cake recipe:

- Take flour
- Add eggs, oil, water
- Mix with spoon for two minutes
- Bake at 350 degrees F for 35 minutes
- Let cool

Piping in tidyverse:

```
flour %>%  
  add(eggs, oil, water) %>%  
  mix(utensil = 'spoon', time = 20) %>%  
  bake(temp = 350, unit = 'F', time = 35) %>%  
  let_cool()
```

Tidyverse functions: data is first argument

- `mutate(data, columns):` `mutate(cars, speed.sq = speed^2)`
add a new column(s) to a data frame
- `select(data, columns):` `select(mtcars, mpg, wt, vs, am)`
selects column(s) from data frame

`mtcars %>%
 select(mpg, wt, vs, am)`

`mtcars %>%
 mutate(wtkg = wt*907.185) %>%
 select(mpg, wt.kg)`

Add a column for weight in kg, then give me only mpg and weight in kilograms columns

Grouping with tidyverse

dplyr functions allow you to perform operations on *groups* of data

`group_by(variables)` to group

`summarise()`, `mutate()`, etc. to operate

Base R equivalent is `aggregate()`

```
mtcars %>%  
  group_by(am, vs) %>%  
  summarise(  
    mean.mpg = mean(mpg),  
    sample.size = n(),  
    se.mpg = sd(mpg) / sqrt(sample.size)  
  )
```

	am	vs	mean.mpg	sample.size	se.mpg
1	0	0	15.0	12	0.801
2	0	1	20.7	7	0.934
3	1	0	19.8	6	1.64
4	1	1	28.4	7	1.80