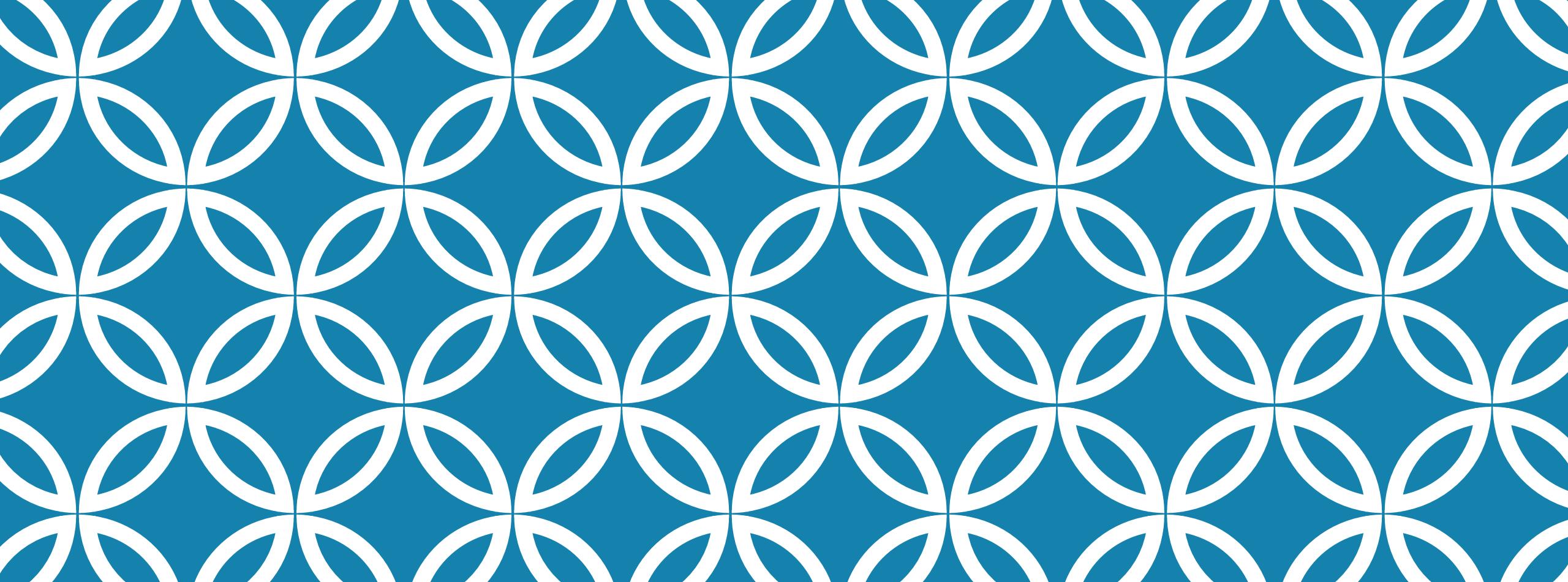


Introduction to R Workshop

August 16, 23 2020



Introduction to R Workshop

Amrom Obstfeld
August 16, 2020



Course Introduction



Goals and Objectives

- Advocate for the use of R as a means of improving reproducibility in clinical data analysis
- Demonstrate how R is used to perform analyses of laboratory operational data
- Establish a basis of understanding in the 'tidy' approach to data analysis within the framework of R



Who are we?

Stephan Kadauke

Assistant Professor of Clinical
Pathology and Laboratory Medicine

University of Pennsylvania Perelman
School of Medicine

Assistant Director of the Cell and
Gene Therapy Laboratory

Children's Hospital of Philadelphia



Daniel Herman

Assistant Professor of Pathology and
Laboratory Medicine

University of Pennsylvania Perelman
School of Medicine

Director, Endocrinology Laboratory

Hospital of the University of
Pennsylvania



Amrom Obstfeld

Assistant Professor of Clinical
Pathology and Laboratory
Medicine

University of Pennsylvania
Perelman School of Medicine

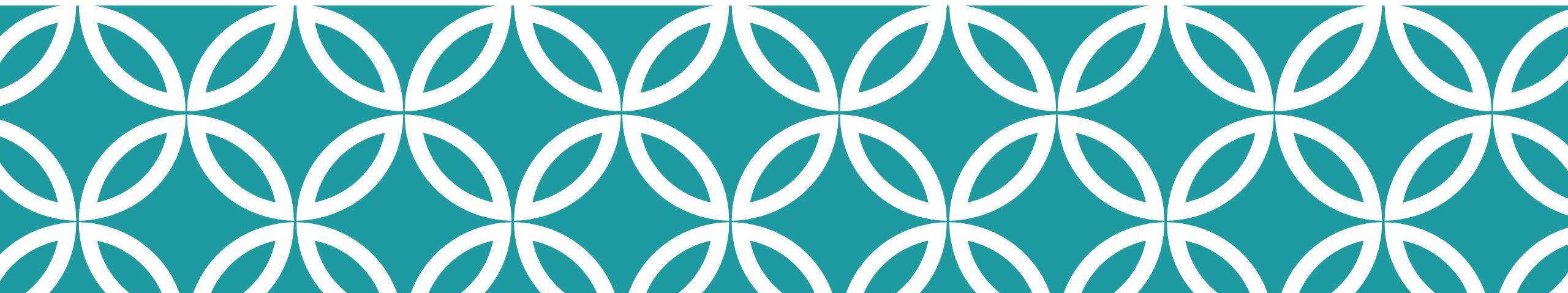
Director of Pathology Informatics

Children's Hospital of
Philadelphia



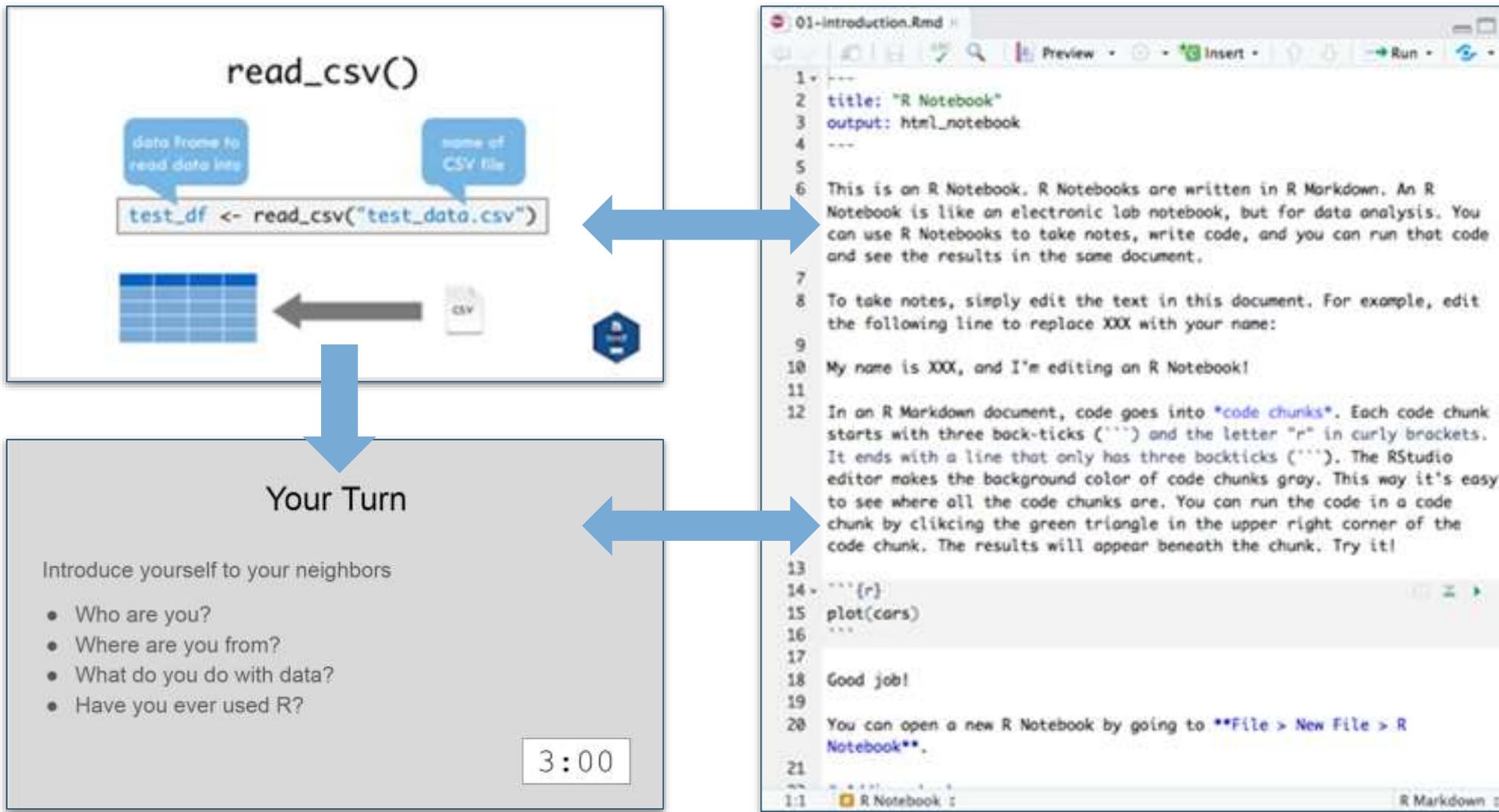


Workshop Workflow



August 16 2020	Session	Instructor
9:00 am - 9:30 am	Instructor Introductions, Introduction to technology	Amrom Obstfeld
9:30 am - 10:15 am	Introduction to R and RStudio	Amrom Obstfeld
10:30 pm - 11:15 am	Reproducible Reporting	Amrom Obstfeld
11:30 am - 1:00 am	Data Visualization	Stephan Kadauke
August 23 2020		
9:00 am - 10:30 pm	Data Transformation	Amrom Obstfeld
10:45 am - 12:15 pm	Statistical Analysis	Dan Herman
12:30 pm - 1:00 pm	Workshop Close out	Amrom Obstfeld

Sessions



Course Project

ESR Reference Range Study

Background and Objectives

One of our clinicians let us know that a lot of the patients, especially healthy elderly ones, receive ESR results flagged as abnormally high. We decided that our reference ranges are not appropriate for this population.

The ESR is the measurement of the rate at which red blood cells (RBCs) settle in anticoagulated blood. RBCs settle faster when pro-inflammatory factors, such as ferritin, cause RBCs to stick together in form rouleaux. Pro-inflammatory factors tend to be elevated in inflammatory states. Thus, the ESR is a non-specific marker for inflammation.

It is known that a normal ESR value varies with age (higher in younger individuals) and gender (higher in females). However, there are multiple reference ranges in use by different clinical laboratories.

It is questionable whether inspection of a mixed population of ESR values (i.e., from healthy and sick individuals) allows classification of normal vs abnormal values. This may be the case if there is clear separation of normal and abnormal values. In any case, it may be useful to visualize a sample of ESR values and highlight different options for cutoffs to better understand whether changing our reference range might be beneficial.

MGH

Age	Male	Female
<1	0-10 mm/h	0-20 mm/h

ARUP Laboratories

Age	Male	Female
<1	0-10 mm/h	0-20 mm/h

Quest Diagnostics

Age	Male	Female
<1	0-10 mm/h	0-20 mm/h
>1	0-20 mm/h	0-30 mm/h

LabCorp

Age	Male	Female
<1	0-10 mm/h	0-20 mm/h
>1	0-20 mm/h	0-30 mm/h

Bakerman's ABCs of Interpretive Laboratory Data

Age	Male	Female
Median	0.2 mm/h	0.2 mm/h
Normal and children	0-10 mm/h	0-15 mm/h
<1	0-10 mm/h	0-20 mm/h
>1	0-20 mm/h	0-30 mm/h

Data Acquisition

Connect to the IMaH Pathology Database (<http://gigex2.uth.tmc.edu:8080>), and select the table `esrresults_submissions`.

Remove all ESR values from 2017, along with the following columns:

1. `Code`: Patient age at time of collection.
2. `CollectionTime`: Date and time of collection.
3. `PtName`: Patient name.
4. `PtNumber`: Patient MRN (AGH).

Data		Actions	
	Value	Count	Operations
CollAge	100	1	Drop
CollAge	47	1	Drop
CollAge	38	1	Drop
CollAge	22	1	Drop
CollAge	39	1	Drop
CollAge	38	1	Drop
CollAge	34	1	Drop
CollAge	41	1	Drop
CollAge	8	1	Drop
CollAge	23	1	Drop
CollAge	27	1	Drop
CollDate	2023-03-01T00:00:00Z	1	Drop
PPhone	1234567890	1	Drop
PNumber	1234567890	1	Drop
PDate	2023-03-01	1	Drop
Result	High	1	Drop
TestOrderName	Test order name	1	Drop
UserGAPFlags	High and Low	1	Drop
Rows	9 rows	9	Drop

The 100% data is stored in the data frame - `acc`. Save the unchanged data frame or disk it (use the database server path below).

Disconnect from the database and clean up.

Data Exploration and Cleaning

The following is a column-by-column inspection and, if necessary, clean up of the data.

Column	Value	Count	Operations
CollAge	100	1	Drop
CollAge	47	1	Drop
CollAge	38	1	Drop
CollAge	22	1	Drop
CollAge	39	1	Drop
CollAge	38	1	Drop
CollAge	34	1	Drop
CollAge	41	1	Drop
CollAge	8	1	Drop
CollAge	23	1	Drop
CollAge	27	1	Drop
CollDate	2023-03-01T00:00:00Z	1	Drop
PPhone	1234567890	1	Drop
PNumber	1234567890	1	Drop
PDate	2023-03-01	1	Drop
Result	High	1	Drop
TestOrderName	Test order name	1	Drop
UserGAPFlags	High and Low	1	Drop
Rows	9 rows	9	Drop

1-10 of 48,573 rows (1-9 of 9 columns)

Previous: 1 2 3 4 5 6 7 8 100 Next

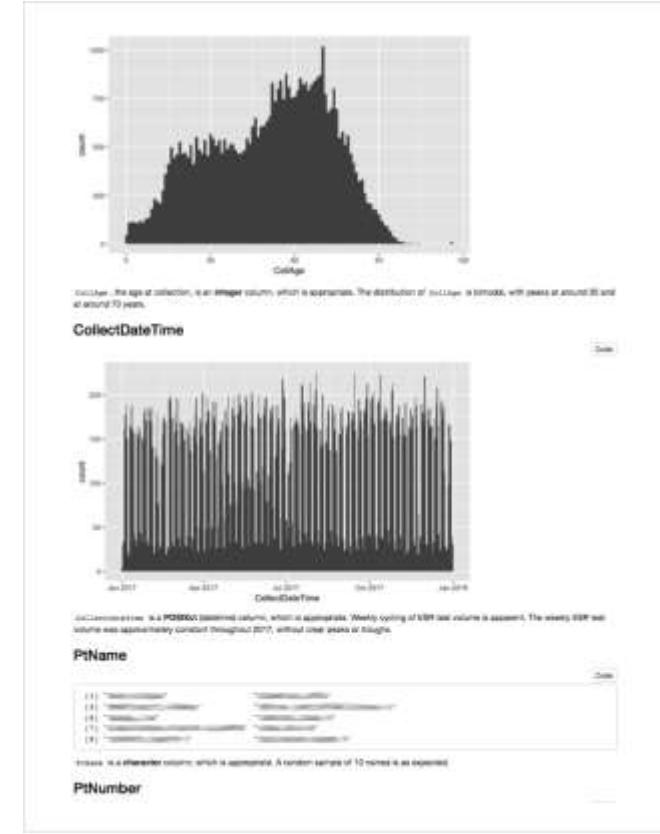
Missing values

Count the number of nulls (missing values) in each column of `acc`.

Column	Nulls	Operations
CollAge	100	Drop
CollDate	0	Drop
PPhone	0	Drop
PNumber	0	Drop
PDate	0	Drop
Result	0	Drop
TestOrderName	0	Drop
UserGAPFlags	0	Drop
Rows	0	Drop

None of the columns of `acc` contain any missing values.

CollAge

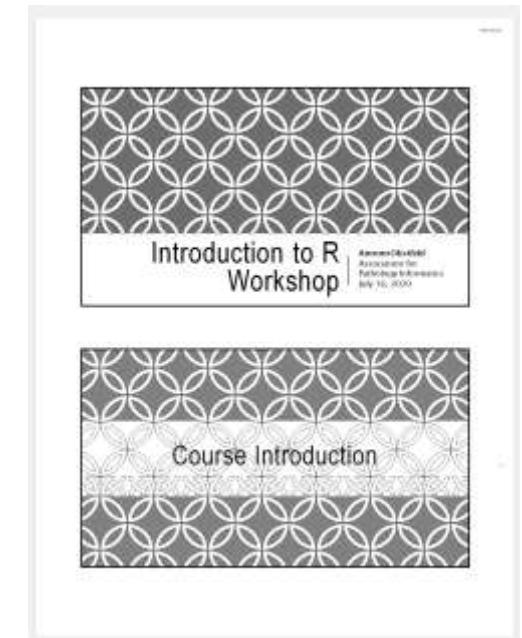
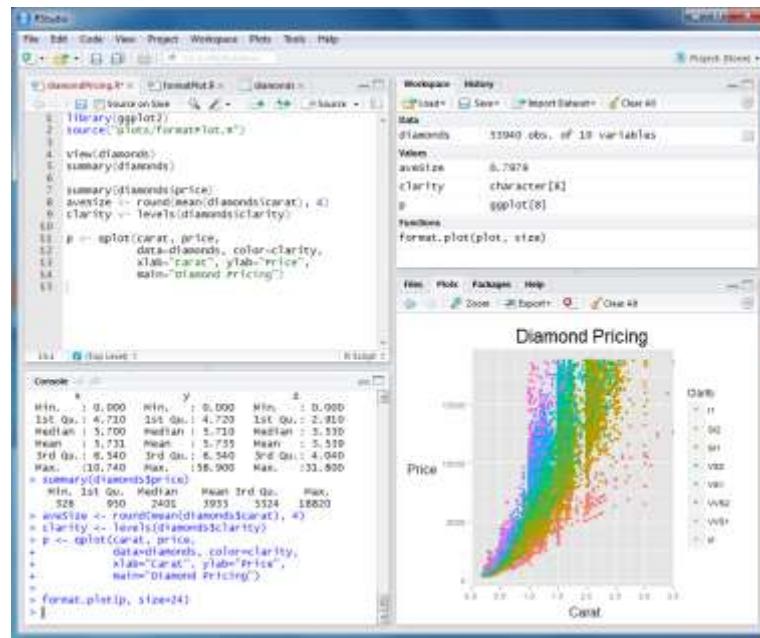


Hackathons

- Participants work with experienced R users on the course project
- Helps troubleshoot problems
- Practice collaborating on analytics projects

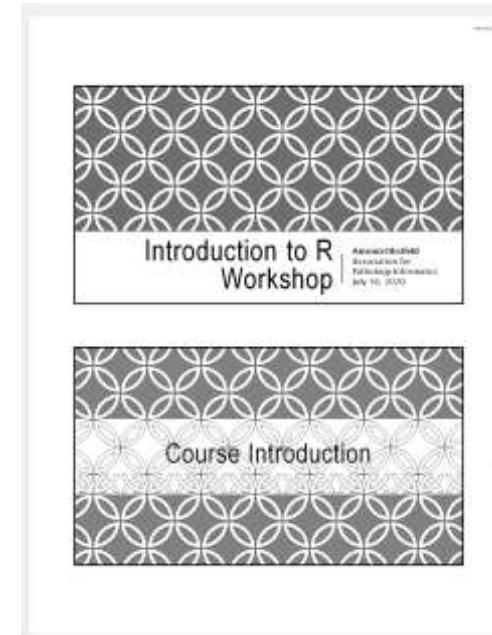


Your Setup

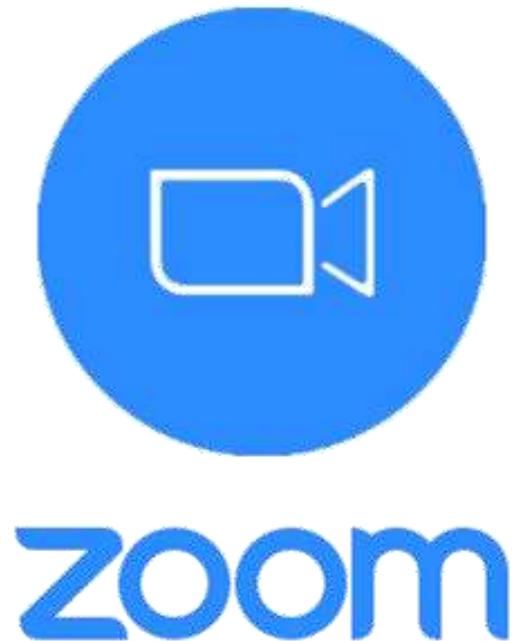


Workshop Coursebook

- Print out of all slides
- Appendix
 - Cheat sheets
 - Useful resources



Using Zoom



- Participants muted
- Chat window
- Non-verbal feedback
- Breakout sessions

Amrom



Using Zoom



zoom

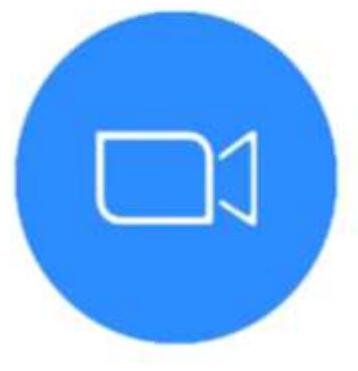
- Participants muted
- Chat window
- Non-verbal feedback
- Breakout sessions

Unmute
Start Video
Participants
Chat
Share Screen
Record
Reactions

Leave

Amrom

Using Zoom



zoom

- Participants muted
- Chat window
- Non-verbal feedback
- Breakout sessions

Participants
2

Chat

Share Screen

Record

Reactions

Start Video

Leave

Participants (2)

Profile	Name	Options
A	Amrom (Me)	...
NS	Nova Smith (Host)	...

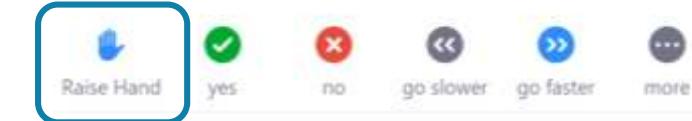
Raise Hand Yes No Go slower Go faster More

Invite Mute Me

Chat

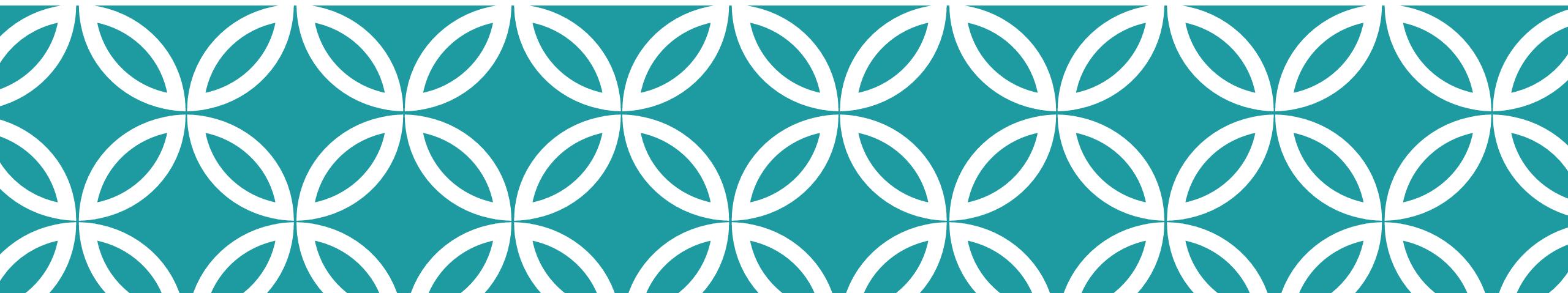
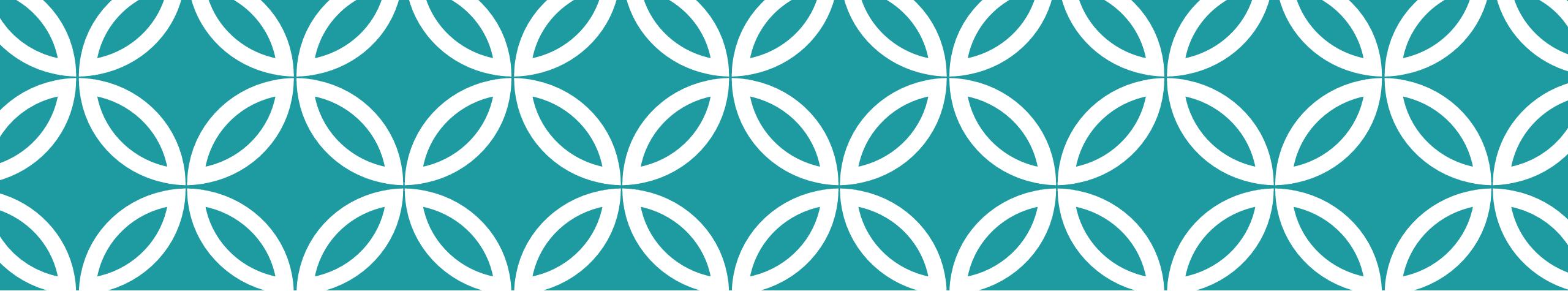
To: Everyone

Type message here...

 Amrom (Me)   Nova Smith (Host)  

Getting Help

- During presentation – Raise hand, instructor will DM
- Break out sessions – Instructor available, unmuted



Who are you?

Your Turn

Introduce yourself!

Who are you?

Where are you from?

Why are you here?

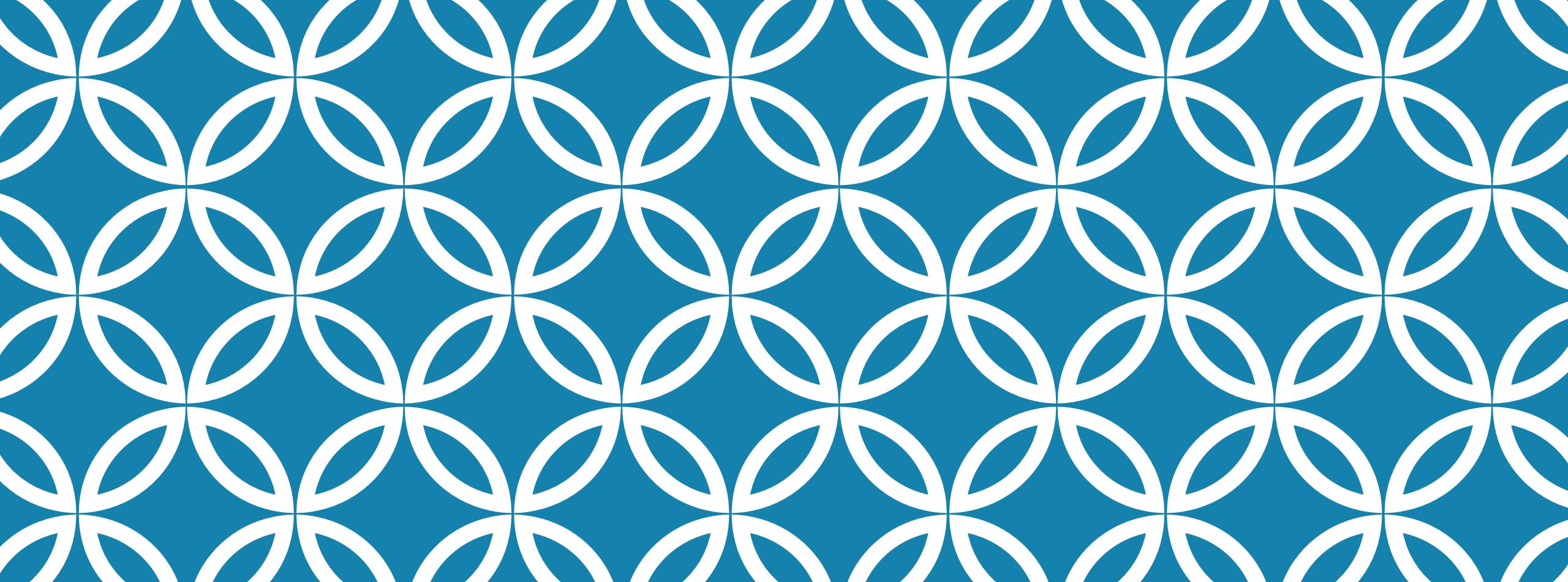
Have you ever used R?

When you need more help

- The Internet (Stack Overflow:
<https://stackoverflow.com/>)
- Work Aids (RStudio Cheat Sheets:
<https://www.rstudio.com/resources/cheatsheets/>)
- A Good Book (R for Data Science:
<http://r4ds.had.co.nz/>)

Final Tips

- The best way to learn to code is by doing
- Practice is key!



Introduction to R and R Studio

Session 1
August 16, 2020
Amrom Obstfeld

August 16 2020

Session

Instructor

9:00 am - 9:30 am Instructor Introductions, Introduction to technology Amrom Obstfeld

9:30 am - 10:15 am Introduction to R and RStudio Amrom Obstfeld

10:30 pm - 11:15 am Reproducible Reporting Amrom Obstfeld

11:30 am - 1:00 am Data Visualization Stephan Kadauke

August 23 2020

9:00 am - 10:30 pm Data Transformation Amrom Obstfeld

10:45 am - 12:15 pm Statistical Analysis Dan Herman

12:30 pm - 1:00 pm Workshop Close out Amrom Obstfeld

Lesson Goals

1. Get oriented to R and RStudio
2. Learn some fundamentals of coding

Lesson Objectives

1. Log in and tour RStudio Cloud
2. Execute code at the console
3. Define and use functions
4. Define and create objects in the environment
5. Load data into R and interact with a dataframe

Getting Oriented to R

What is R?

- R is a statistical programming language.
- Using R you can load, analyze, and visualize data.
- R also provides an environment in which we can conduct reproducible data analysis.
 - Documented
 - Revisable
 - Shareable



RStudio: The Portal to R

- RStudio is an integrated development environment (IDE)
- Using RStudio we can interact with the R programming language to:
 - Write and execute code interactively
 - View data
 - Debug and fix errors
 - Author our code



RStudio: In the Cloud... In Your Home



- RStudio Cloud: An online hosted version of RStudio that we will use for these course sessions
- RStudio Desktop: A locally installed version of RStudio that you will use when you get home to continue your learning



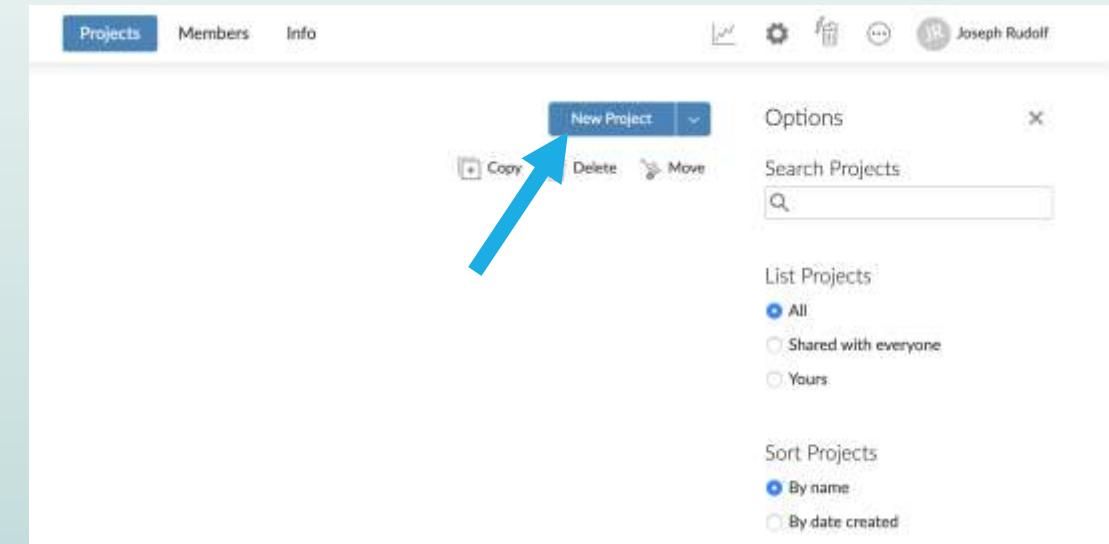
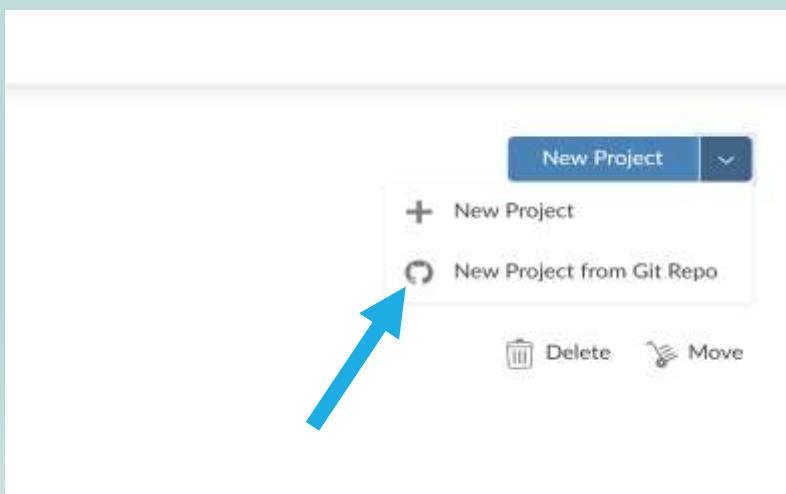
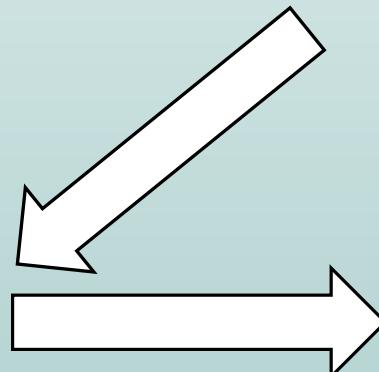
Note: Use Rstudio Cloud only for this course. Do not upload protected health information to the cloud!

Your Turn

Navigate to:
rstudio.cloud



Register for an account



Type in this link at the prompt:
<https://tinyurl.com/pennres2020>



Spaces

[Your Workspace](#)[AACC 2019 Introduction to R](#)[API R Workshop 2020](#)[New Space](#)

Learn

[Guide](#)[What's New](#)[Primers](#)[Cheat Sheets](#)[Feedback and Questions](#)

Info

[Plans & Pricing](#)[Terms and Conditions](#)[System Status](#)[File](#) [Edit](#) [Code](#) [View](#) [Plots](#) [Session](#) [Build](#) [Debug](#) [Profile](#) [Tools](#) [Help](#)[+ New File](#) [Open](#) [Save](#) [Save As](#) [Import](#) [Export](#) [Go to file/function](#) [Addins](#)

R 4.0.0

[Console](#) [Terminal](#) [Jobs](#)

/cloud/project/

R version 4.0.0 (2020-04-24) -- "Arbor Day"
Copyright (C) 2020 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> |

[Environment](#) [History](#) [Connections](#) [Git](#)[Import Dataset](#)

Global Environment

Environment is empty

[Files](#) [Plots](#) [Packages](#) [Help](#) [Viewer](#)[New Folder](#) [Upload](#) [Delete](#) [Rename](#) [More](#)

Cloud > project

Name	Size	Modified
..		
.gitignore	621 B	Jun 23, 2020, 9:21 PM
.Rhistory	0 B	Jul 13, 2020, 1:26 PM
.Rprofile	88 B	Jun 23, 2020, 9:25 PM
03 - Visualize.Rmd	3 KB	Jul 10, 2020, 8:46 AM
04 - Transform.Rmd	4.8 KB	Jul 13, 2020, 12:08 PM
05 - Stats.Rmd	5.8 KB	Jul 10, 2020, 8:46 AM
06 - Advanced Reporting.Rmd	871 B	Jul 13, 2020, 7:08 AM
coursepack		
data		
LICENSE	1 KB	Jun 23, 2020, 9:21 PM
presentations		
project.Rproj	205 B	Jul 13, 2020, 1:26 PM
README.md	6.9 KB	Jul 12, 2020, 3:42 PM

The image shows the RStudio interface with several panes:

- Editor**: The top-left pane contains an untitled R script file. The code area has lines 1 and 2 visible. The toolbar includes icons for saving, running, and source code.
- Console**: The bottom-left pane shows the current working directory as `~/Rconnect/meditech/`.
- MISC**: The bottom-right pane displays the RStudio Connect environment. It shows the "Global Environment" tab with the message "Environment is empty". Below it, the "Files" tab shows a directory structure:
 - ..
 - rstudio-connect.chop.edu



The Basics of Coding

The Basics of Coding: Calculation

- R is a calculator!

```
> 2 + 3 + 2  
[1] 7  
>  
>  
> 4 * 20  
[1] 80  
>  
>  
> 6 ^ 8  
[1] 1679616  
>
```

enter/return
to execute
code

answer
returned here

Your Turn #1

Place your cursor at the console and click to enter the console.

Complete the following calculation:

- For the date 12-29-1974
- Take the four digit year
- Subtract the month then multiply by the day

What did you get?

- A four digit number? A five digit number?

```
> 1974 - 12 * 29  
[1] 1626  
>  
>  
> (1974 - 12) * 29  
[1] 56898
```

- Order of operations matters!

The Basis of Coding: 3 Players

1. Functions
2. Arguments
3. Objects

The Basics of Coding: Functions

- Code that extends our reach beyond the basic operators

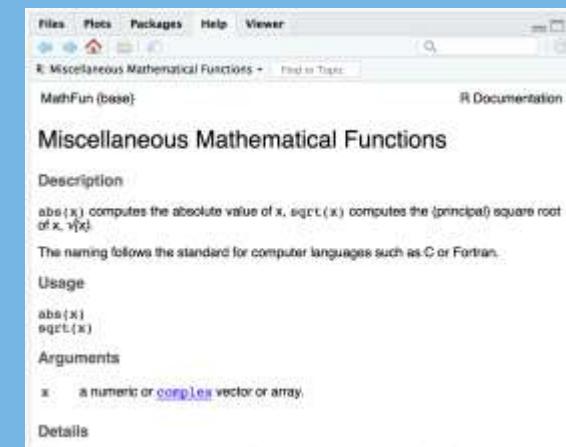
```
> abs(-77)  
[1] 77  
>
```

function
(does stuff)

abs(-77)

- What if I don't know what a function does?

```
> ?abs()  
>
```



The Basics of Coding: Arguments

- The input that defines what the function should do

```
> abs(-77)
```

```
[1] 77
```

```
>
```

argument
(input)

```
abs(-77)
```

The Basics of Coding: Objects

- Objects are the container for your output

object
(stores
output)

```
my_abs <- abs(-77)
```

Checking the contents of an object

- Entering the object name at the console allows us to output the contents of an object.

```
> abs(-77)
[1] 77
> my_abs <- abs(-77)
> my_abs
[1] 77
> |
```

Checking the contents of an object

- The environment tab shows us the objects we have created.

A screenshot of the RStudio interface, specifically the Environment tab. The tab bar includes tabs for Environment, History, Connections, Git, and Tutorial. The Environment tab is active, indicated by a thicker border. Below the tabs, there are several icons: a folder with a green arrow, a blue floppy disk, a grid with a green arrow, and a broom. Next to these icons is a dropdown menu labeled "Import Dataset". Below the toolbar, there is a section titled "Global Environment" with a dropdown arrow. The main area is titled "Values" and lists one item: "my_abs". At the bottom right of the main area, the number "77" is displayed, likely indicating the total number of objects in the environment or a specific identifier. The overall background is light gray, and the interface has a clean, modern look.

Bending objects to your will

- Once we have created an object we can start to interact with it.
- This includes passing our objects to other functions... Whoa!

```
> log(my_abs, 2)
[1] 6.266787
>
```

Knowledge Check

Consider this code:.

```
mean_age <- mean(age, na.rm=TRUE)
```

Which is the function?

Which is the argument?

Which is the objects?

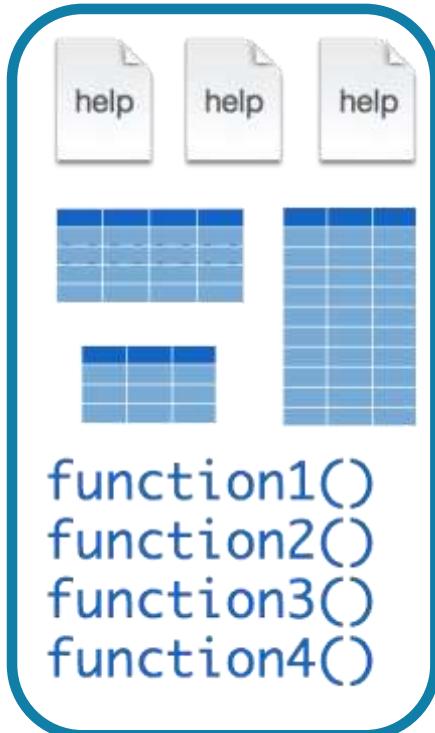
The Basics of Coding: Packages

- A package is a collection of functions.
- Packages extend the capabilities of the base R programming language.
- The **tidyverse** includes functions for reading data into the R environment, cleaning and manipulating data, and plotting our results.



A Word About Packages

tidyverse



function1()
function2()
function3()
function4()

1

`install.packages("tidyverse")`

Downloads files to computer
1 x per computer

2

`library("tidyverse")`

Loads package
1 x per R Session

Your Turn #2

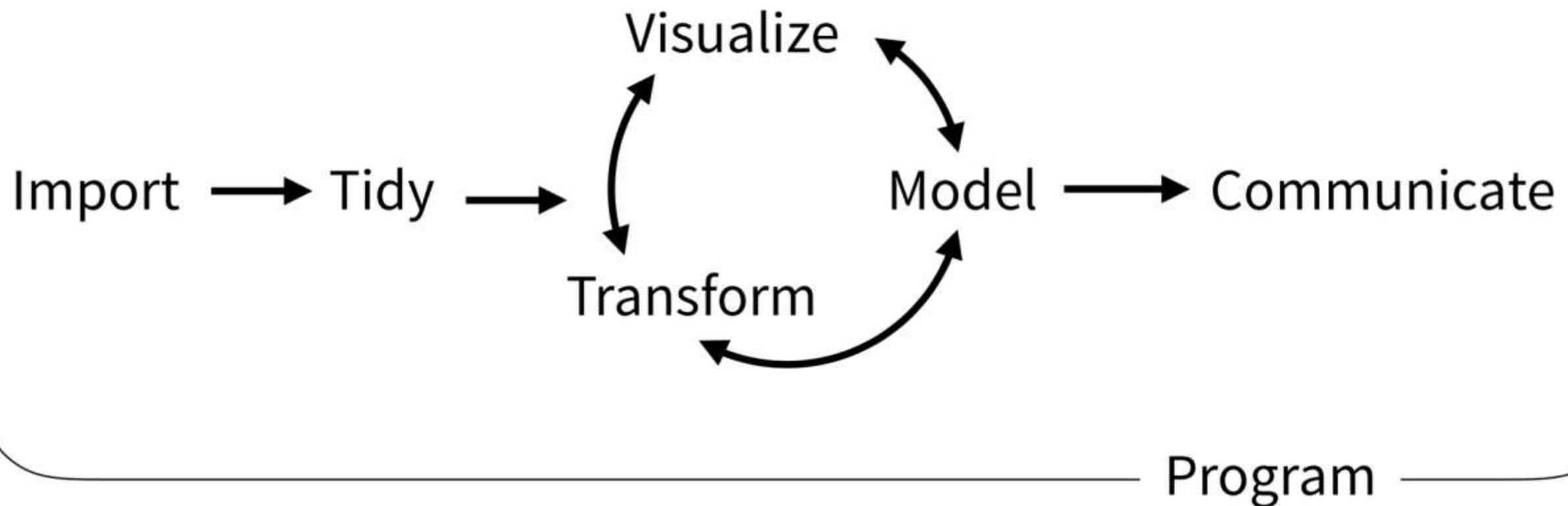
Run the following in the console:

```
install.packages("tidyverse")
```

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

The Data Analysis Process

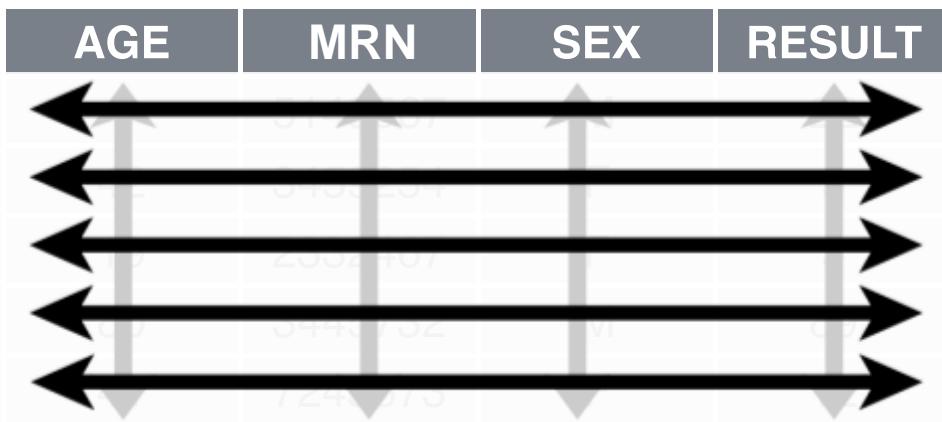
Typical Data Analysis Pipeline



What is a “Tidy” Data Frame

A data set is **tidy** if:

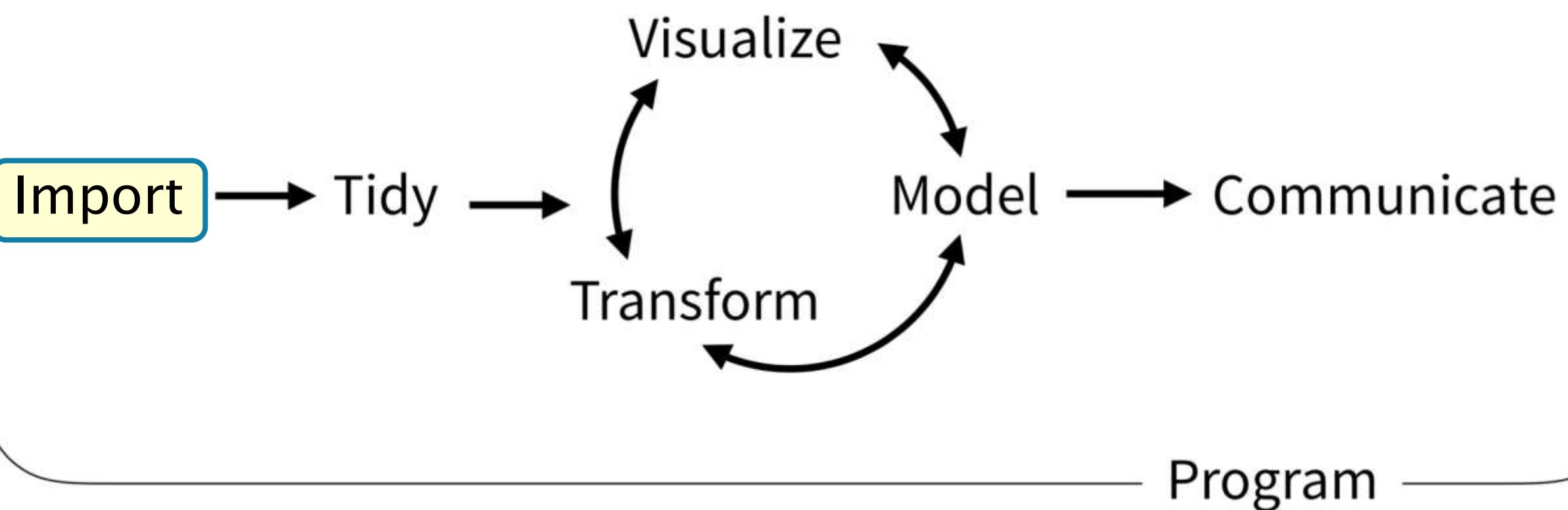
AGE	MRN	SEX	RESULT
1	100001	M	POSITIVE
2	100002	F	NEGATIVE
3	100003	M	POSITIVE
4	100004	F	NEGATIVE



1. Each **variable** is in its own **column**
2. Each **observation** is in its own **row**
3. Each **value** is in its own **cell**

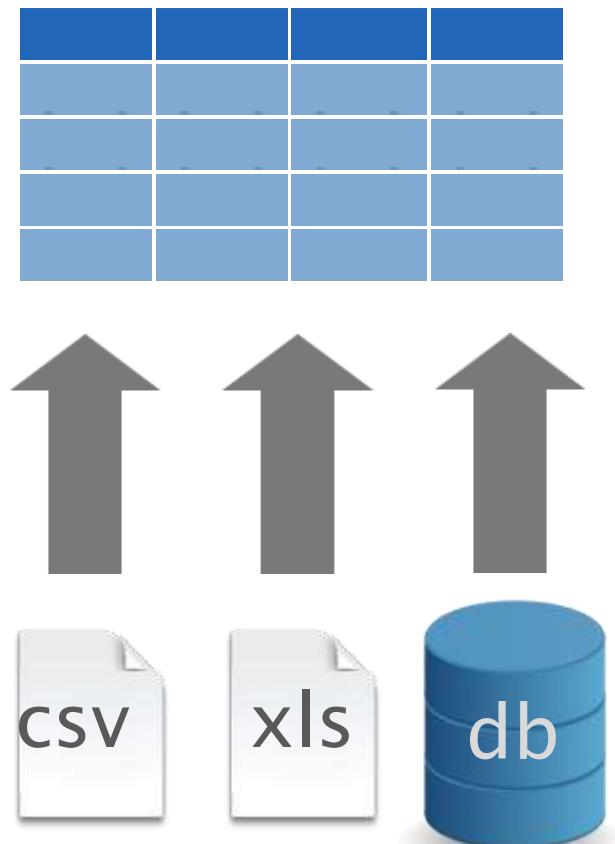
Importing Data

Typical Data Analysis Pipeline



Dataframes: Beyond the Vector

- Dataframe is the term for a table
- Dataframes are composed:
Columns (Variables)
Rows (Observations)
- Dataframes are objects and can be
acted on like other objects



plain text
("flat") file

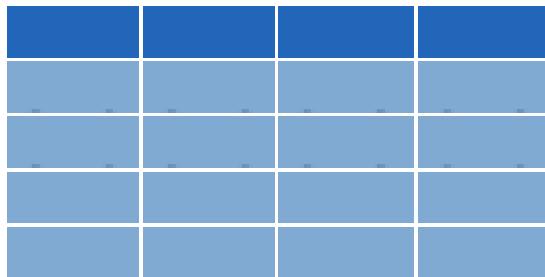
header
row

Name	MRN	DOB
Santa Claus	12345	1/1/01
Roger Rabbit	67890	12/12/69
Kermit the Frog	24680	2/2/22

rectangular
structure

Loading Data to Create a Dataframe

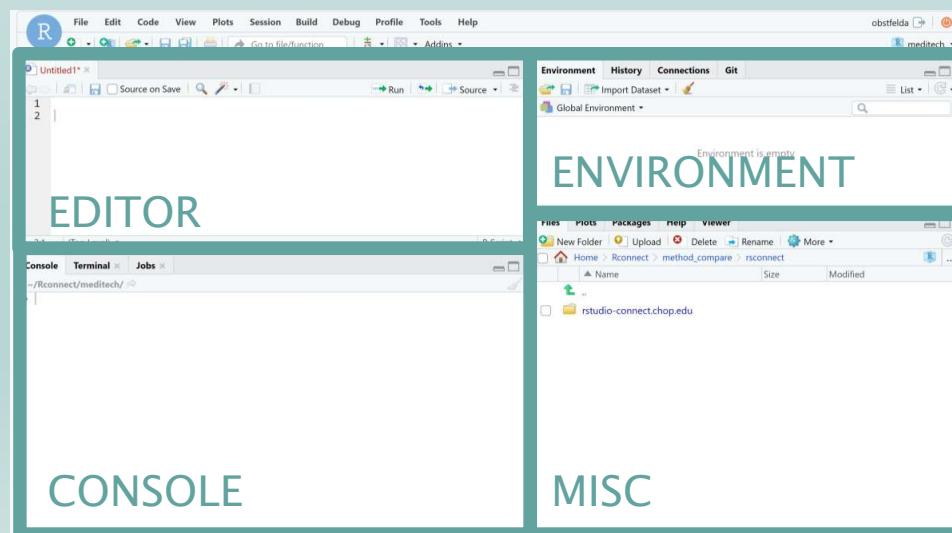
```
data_frame <- read_csv("file_name")
```



Memory Check

After reading in data using code such as this,
where will you the data appear?

```
data_frame <- read_csv("file_name")
```



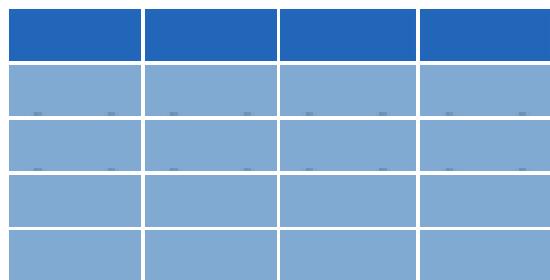
`read_csv()`

data frame
to read data
into

name of
CSV file

```
covid_testing <- read_csv("data/covid_testing.csv")
```

`covid_testing`



`covid_testing.csv`



Your Turn #3

Configure environment and load the Covid Testing CSV:

Load the tidyverse library using `library(tidyverse)`

Use the `read_csv()` function to load the data

-File_name argument: “`data/covid_testing.csv`”

-Object name: `covid_testing`

What's in a name?

- Capitalization matters

covid_testing

\neq

Covid_testing

\neq

COVID_TESTIN

G

- Strive for names that are concise and meaningful (not easy!)

Bad

p

Still not great

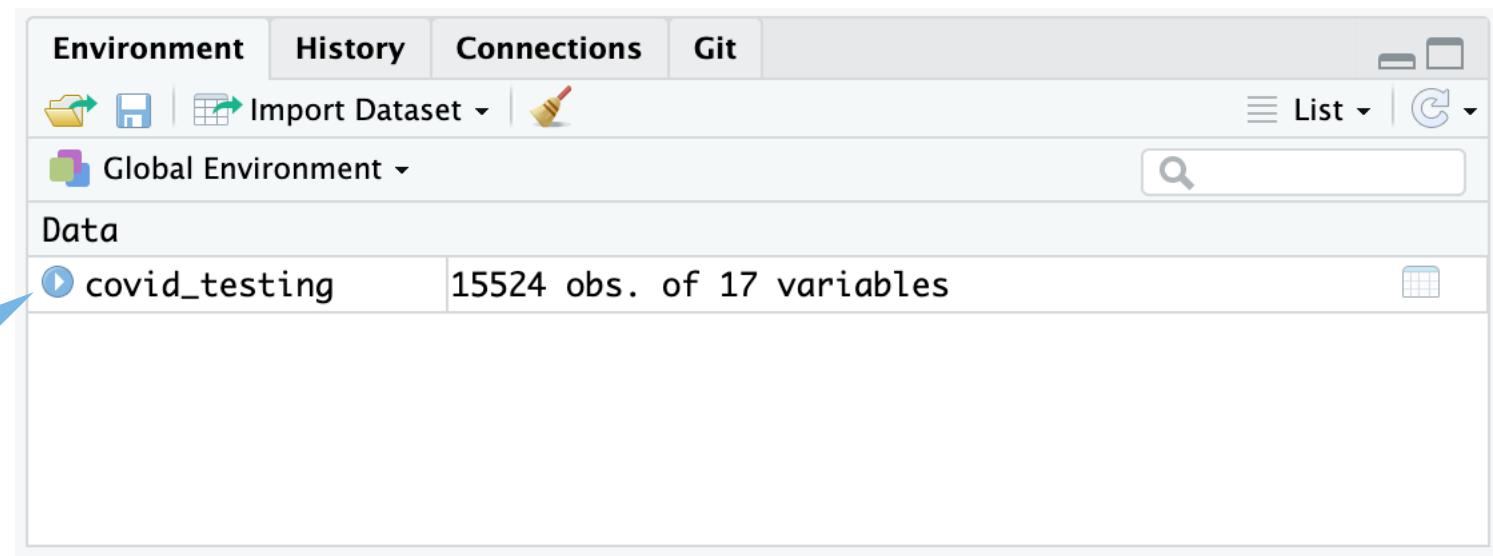
name

Good

patient_name

Viewing the Contents of a Dataframe

single click to explore the data



Viewing the Contents of a Dataframe

	mrn	first_name	last_name	gender	pan_day	test_id	clinic_name	result
1	5001412	jhezane	westerling	female	4	covid	inpatient ward	positive
2	5000533	penny	targaryen	female	7	covid	clinical lab	negative
3	5009134	grunt	rivers	male	7	covid	clinical lab	negative
4	5008518	melisandre	swyft	female	8	covid	clinical lab	negative
5	5008967	rolley	karstark	male	8	covid	emergency dept	positive
6	5011048	megga	karstark	female	8	covid	oncology day hosp	negative
7	5000663	ithoke	targaryen	male	9	covid	clinical lab	negative
8	5002158	ravella	frey	female	9	covid	emergency dept	negative
9	5003794	styr	tyrell	male	9	covid	clinical lab	negative
10	5004706	wvnafrvd	seaworth	male	9	covid	clinical lab	negative

15,524
Observations
(Rows)

17 Attributes
(Columns)

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 = lappend)
```

File with arbitrary delimiter

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

CSV for excel

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

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 = lappend)
```



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())
```

Comma Delimited Files

```
read_csv("file.csv")  
# To make file.csv run:  
write_file(x = "abc\n1,2,3\n4,5,NA", path = "file.csv")
```

Semi-colon Delimited Files

```
read_csv2("file2.csv")  
# write_file(x = "abc;c1,2;3\n4;5;NA", path = "file2.csv")
```

Files with Any Delimiter

```
read_delim("file.txt", delim = "|")  
# write_file(x = "abc|c1|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

Example file
write_file("a,b,c\n1,2,3\n4,5,NA", "file.csv")
file.csv

No header
read_csv(l, col_names = FALSE)
l = file.csv

Provide header
read_csv(l, col_names = c("x", "y", "z"))
l = file.csv

Missing Values
read_csv(l, na = c("1", "0"))
l = file.csv

Read Non-Tabular Data

Read a file into a single string

```
read_file(file, locale = default_locale())
```

Read each line into its own string

```
read_lines(file, skip = 0, n_max = -1, 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(l)

Read each line into a raw vector
read_lines_raw(file, skip = 0, n_max = -1,
 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:  
# col_integer()  
# age = col_integer(),  
# sex = col_character(),  
# earn = col_double()  
# ...  
# age is a double (numerical)  
# sex is a character
```

1. Use `problems()` to diagnose problems

```
#> read_csv("file.csv"); problems()
```

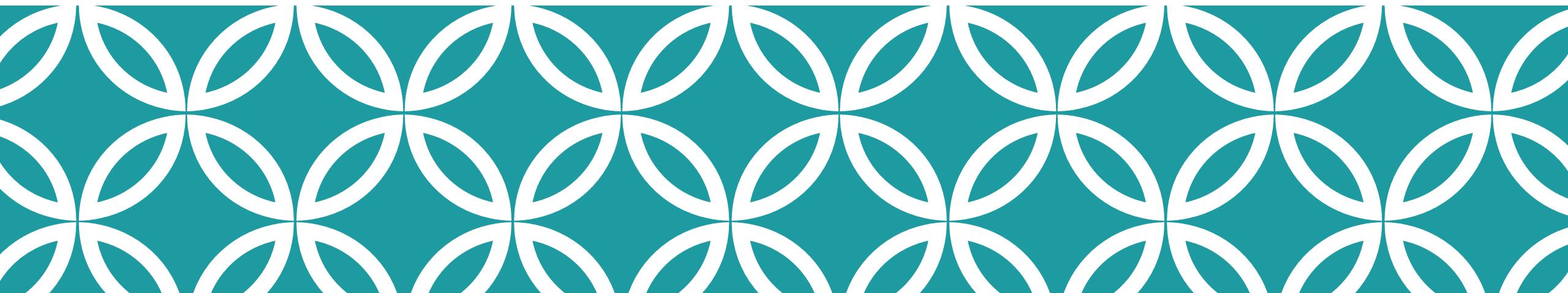
2. Use a `col_` function to guide parsing

- `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()`
- #> read_csv("file.csv", col_types = col_id(
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()`
- #> parse_number(xSA)

RStudio® is a trademark of RStudio, Inc. - CC BY-SA. RStudio | info@rstudio.com | 444 4th Street | RStudio.com | Learn more with tidyverse.org | RStudio 1.1.2 | tidyverse 1.1.2 | tidy 0.8.0 | Updated: 2017-01



What Else?



Import Excel files
(.xls, .xlsx)

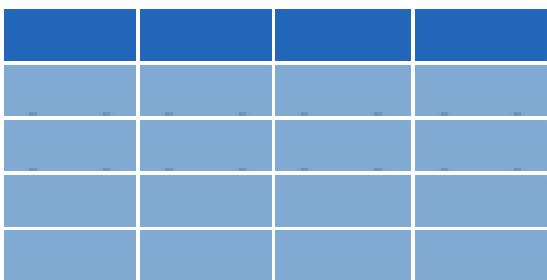
```
library(readxl)
```

read_excel()

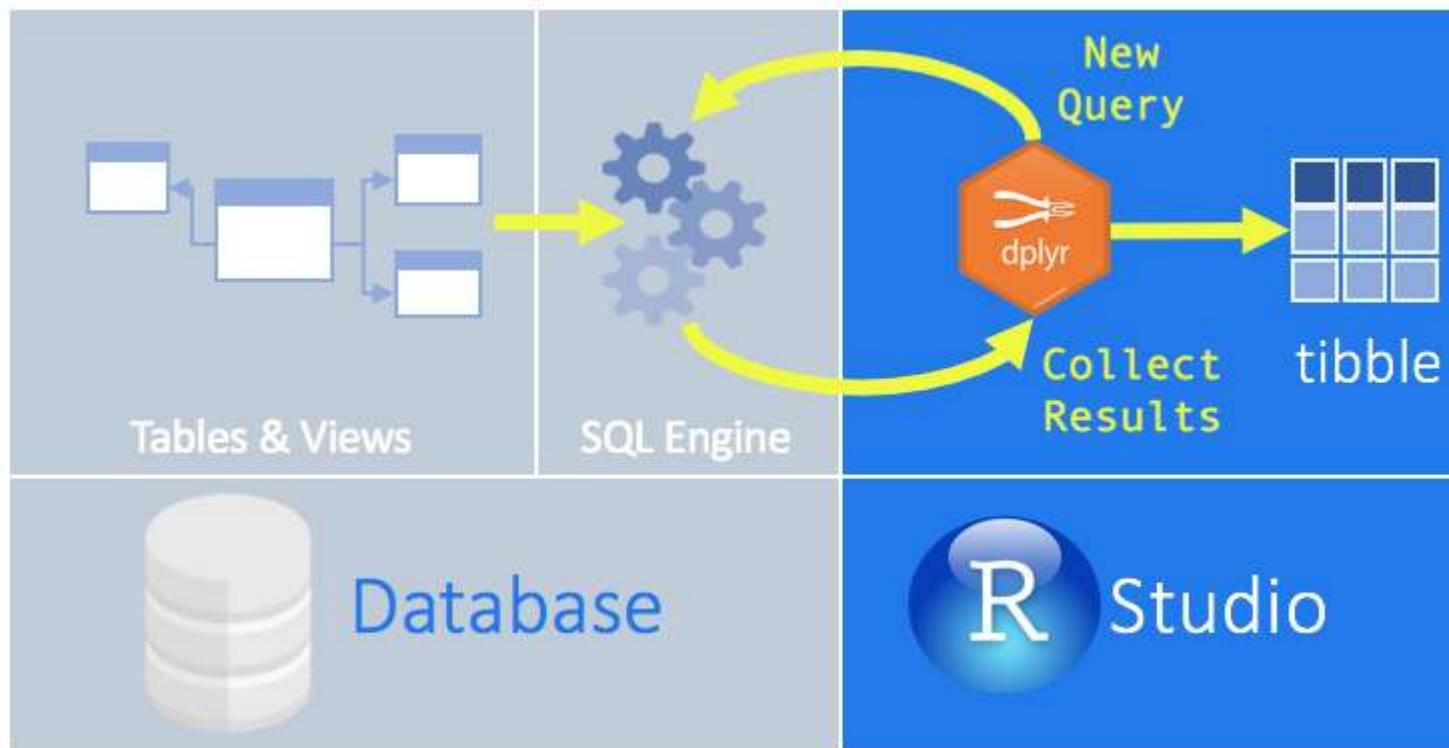
data frame
to read data
into

name of
Excel file

```
data_frame <- read_excel(file_name)
```



Read Directly From Database





Reads SPSS, Stata, and
SAS files



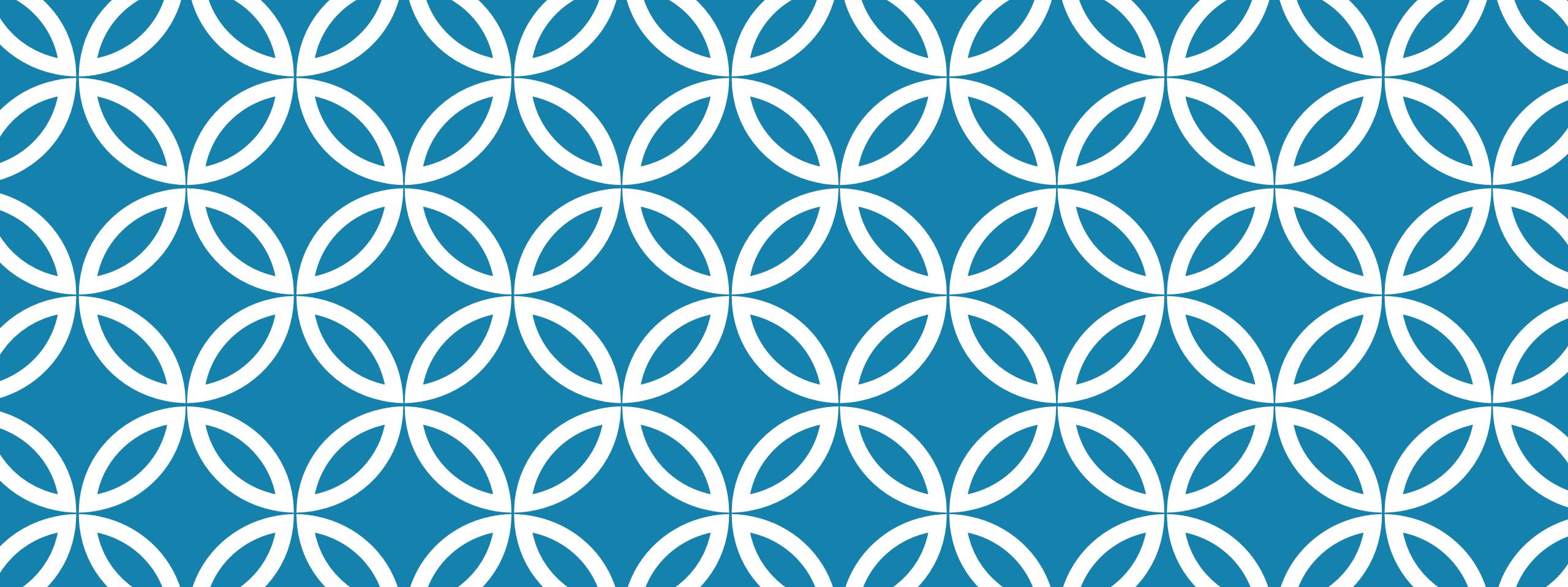
Format	Typical Extension	Import Package	Export Package	Installed by Default
Comma-separated data	.csv	data.table	data.table	Yes
Pipe-separated data	.psv	data.table	data.table	Yes
Tab-separated data	.tsv	data.table	data.table	Yes
CSVY (CSV + YAML metadata header)	.csvy	data.table	data.table	Yes
SAS	.sas7bdat	haven	haven	Yes
SPSS	.sav	haven	haven	Yes
Stata	.dta	haven	haven	Yes
SAS XPORT	.xpt	haven	haven	Yes
SPSS Portable	.por	haven		Yes
Excel	.xls	readxl		Yes
Excel	.xlsx	readxl	openxlsx	Yes
R syntax	.R	base	base	Yes
Saved R objects	.RData, .rda	base	base	Yes
Serialized R objects	.rds	base	base	Yes
JSON	.json	jsonlite	jsonlite	No

Lesson Goals

1. Get oriented to R and RStudio
2. Learn some fundamentals of coding

Lesson Objectives

1. Log in and tour RStudio Cloud
2. Execute code at the console
3. Define and use functions
4. Define and create objects in the environment
5. Load data into R and interact with a dataframe



Reproducible Reporting

Session 2
August 16, 2020
Amrom Obstfeld

August 16 2020	Session	Instructor
9:00 am - 9:30 am	Instructor Introductions, Introduction to technology	Amrom Obstfeld
9:30 am - 10:15 am	Introduction to R and RStudio	Amrom Obstfeld
10:30 pm - 11:15 am	Reproducible Reporting	Amrom Obstfeld
11:30 am - 1:00 am	Data Visualization	Stephan Kadauke
August 23 2020		
9:00 am - 10:30 pm	Data Transformation	Amrom Obstfeld
10:45 am - 12:15 pm	Statistical Analysis	Dan Herman
12:30 pm - 1:00 pm	Workshop Close out	Amrom Obstfeld

Goals

1. Understand why reproducible reporting is important
2. Learn to work within R Markdown for reproducible reports

Objectives

1. Create an R Markdown document and generate different types of output files
2. Practice modifying each component of a R Markdown file

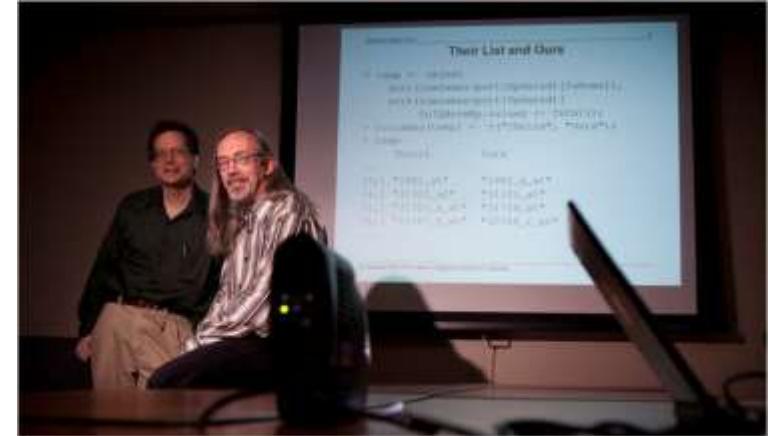
Why Analyze Data with Code?

The Duke Cancer Scandal

- Chemo sensitivity from microarrays
- Errors first, then misconduct
- Clinical trials based on flawed models
- Papers retracted, lawsuits settled



“Common errors are simple,
Simple errors are common”



Theirs

"1881_at"

"31321_at"

"31725_s_at"

"32307_r_at"

...

Ours

"1882_g_at"

"31322_at"

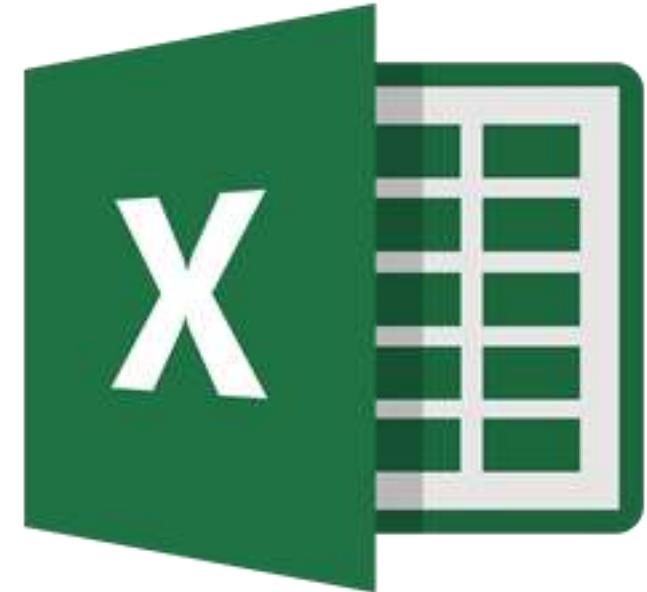
"31726_at"

"32308_r_at"

<https://youtu.be/tN31vZj0LZY>

Point-and-Click Is Not Reproducible

- Interactive tools do not record user actions
- Manual documentation is error-prone
- Manual analyses cannot be repeated on new data sets or shared with collaborators



Computer code can precisely document each step of the analysis

Why YOU Should Do Data Analysis Reproducibly

“Can we redo the analysis with this month’s data?”

“Why do the data in Table 1 not seem to agree with Figure 2?”

“Why did I decide to omit these six samples?”



**YOUR CLOSEST COLLABORATOR IS YOU FROM 6 MONTHS AGO
(BUT YOU DON'T ANSWER E-MAILS)**

Replication vs Reproduction

- ❖ Replication: other people collect new data
 - Scientific gold standard
 - Difficult and time-consuming

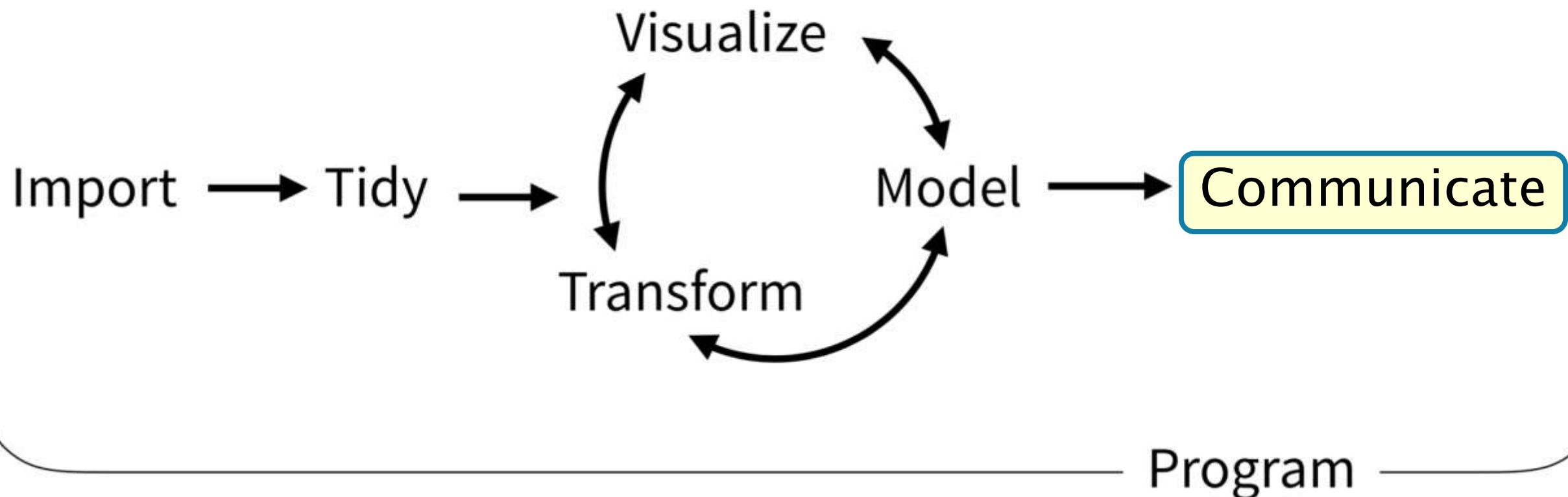
- ❖ Reproduction: other people analyze the same data
 - Does not by itself validate the analysis ...
 - Has been proposed as a minimal standard

Your Turn #1

Spend one minute writing down as many reasons as you can think of for why you might want your own data analyses to be reproducible.

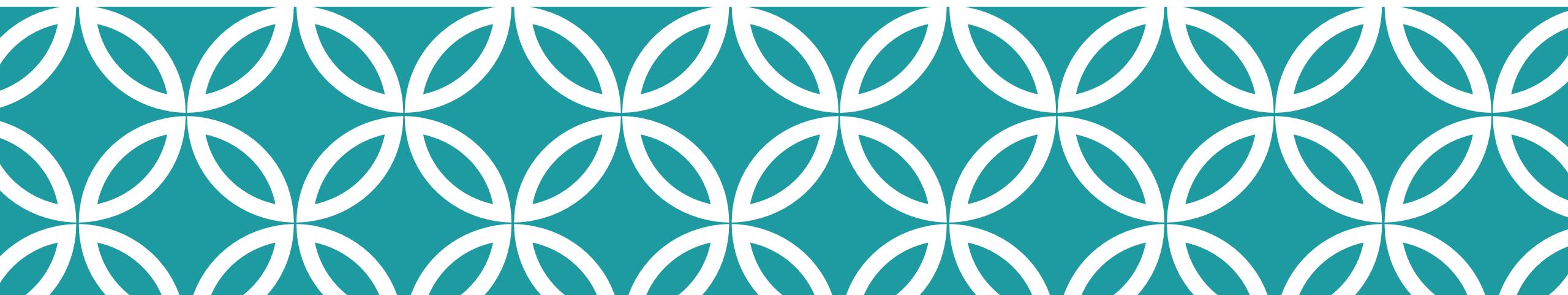


Typical Data Analysis Pipeline





Using R/R Markdown to Support Reproducibility



Using R for Reproducibility

Programming in R (or another language) allows one to reproduce analysis steps exactly or perform same analysis on new data

Better practice is to create documentation about analysis to accompany and explain code

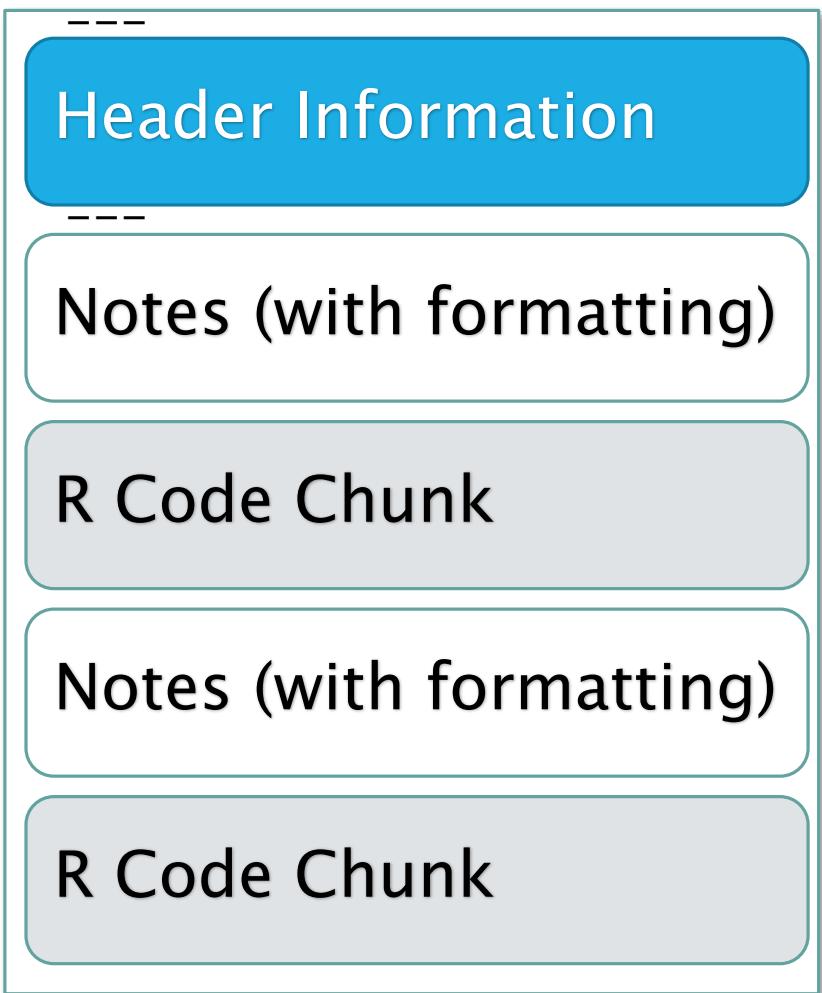
Best practice is include documentation and code in one place

R Markdown for Reproducible Reports

- Code remains attached to its documentation and output
- Text and code compile into a single HTML or PDF document that can be shared
- Reproducible reports can be automated and turned into analytic dashboards

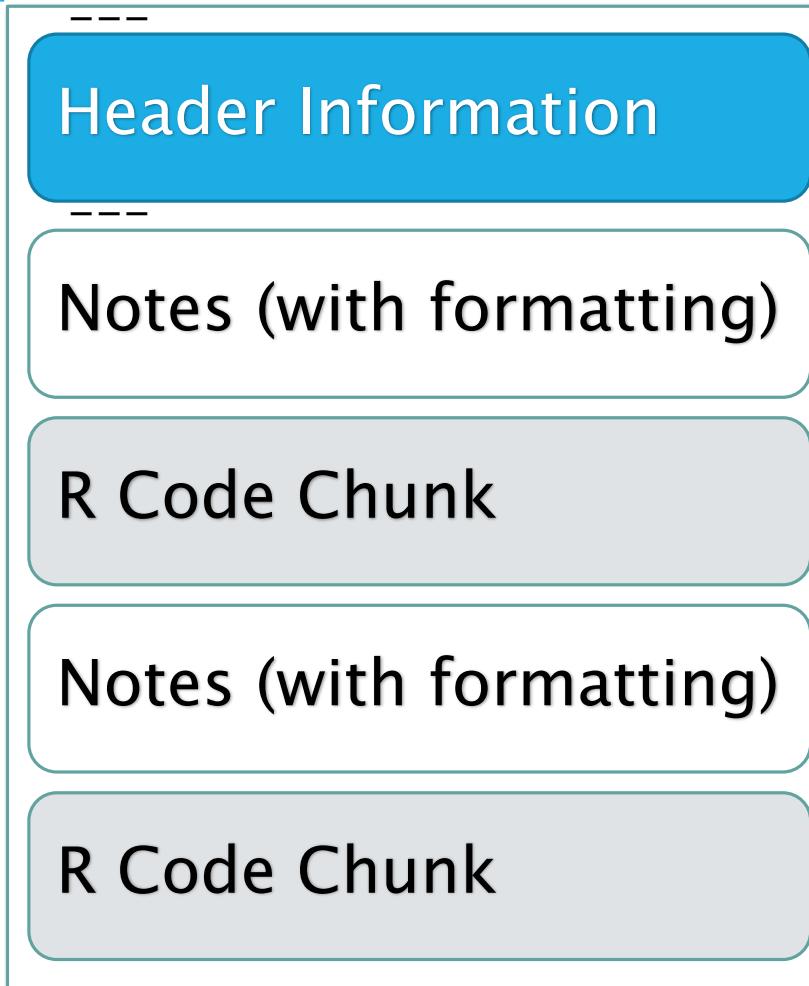


R Markdown Block Diagram



```
Untitled1  
Process the whole document Insert +  
1+ ---  
2+ title: "notes"  
3+ author: "Kieran healy"  
4+ date: "9/13/2017"  
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 form of authoring HTML, PDF, and MS Word documents. For more details 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 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 printing of the R code that generated the plot.  
2.1 notes 2
```

R Markdown Block Diagram



```
1+ ---
2  title: "My Project"
3  author: "Julie"
4  date: "11/21/2017"
5  output: html_document
6  ---
13
14 This is an R Markdown document. Markdown is a simple formatting syntax for authoring 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 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
```

```
# One Hashtag = Large Header
```

```
## Two Hashtags = Smaller Header
```

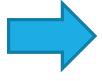
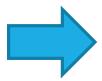
Here is some text.

- * It's easy to make a list.
- * Here's how you style text *cursive* or **bold**.
- * Let's add a [link](https://www.massgeneral.org).

```
```{r}
x <- rnorm(100)
summary(x)
```
```

```
## Including plots
```

```
```{r, echo = FALSE}
hist(x)
```
```



One Hashtag = Large Header

Two Hashtags = Smaller Header

Here is some text.

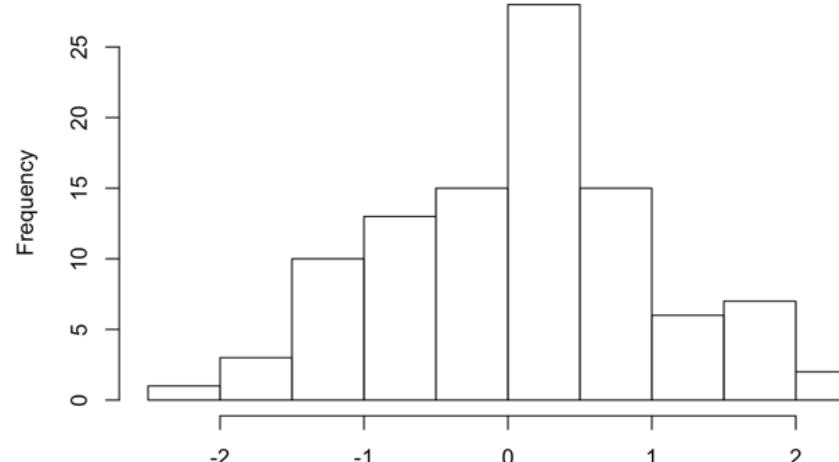
- It's easy to make a list.
- Here's how you style text *cursive* or **bold**.
- Let's add a link.

```
x <- rnorm(100)
summary(x)
```

```
##      Min. 1st Qu. Median      Mean 3rd Qu.      Max.
## -2.0053 -0.6084  0.1125   0.0747  0.6100   2.1437
```

Including plots

Histogram of x





R Markdown

Reference Guide

Learn more about R Markdown at rmarkdown.rstudio.com

Learn more about Interactive Docs at shiny.rstudio.com/articles

Contents:

- 1. Markdown Syntax**
2. Knitr chunk options
3. Pandoc options

Syntax

Plain text

End a line with two spaces
to start a new paragraph.

italics and _italics_

bold and __bold__

superscript^{^2}

~~strikethrough~~

[link] (www.rstudio.com)

Header 1

Header 2

Header 3

Header 4

Becomes

Plain text

End a line with two spaces to start a new paragraph.

italics and *italics*

bold and **bold**

superscript²

strikethrough

link

Header 1

Header 2

Header 3

Code chunks

Open/close
with 3
backticks

Language

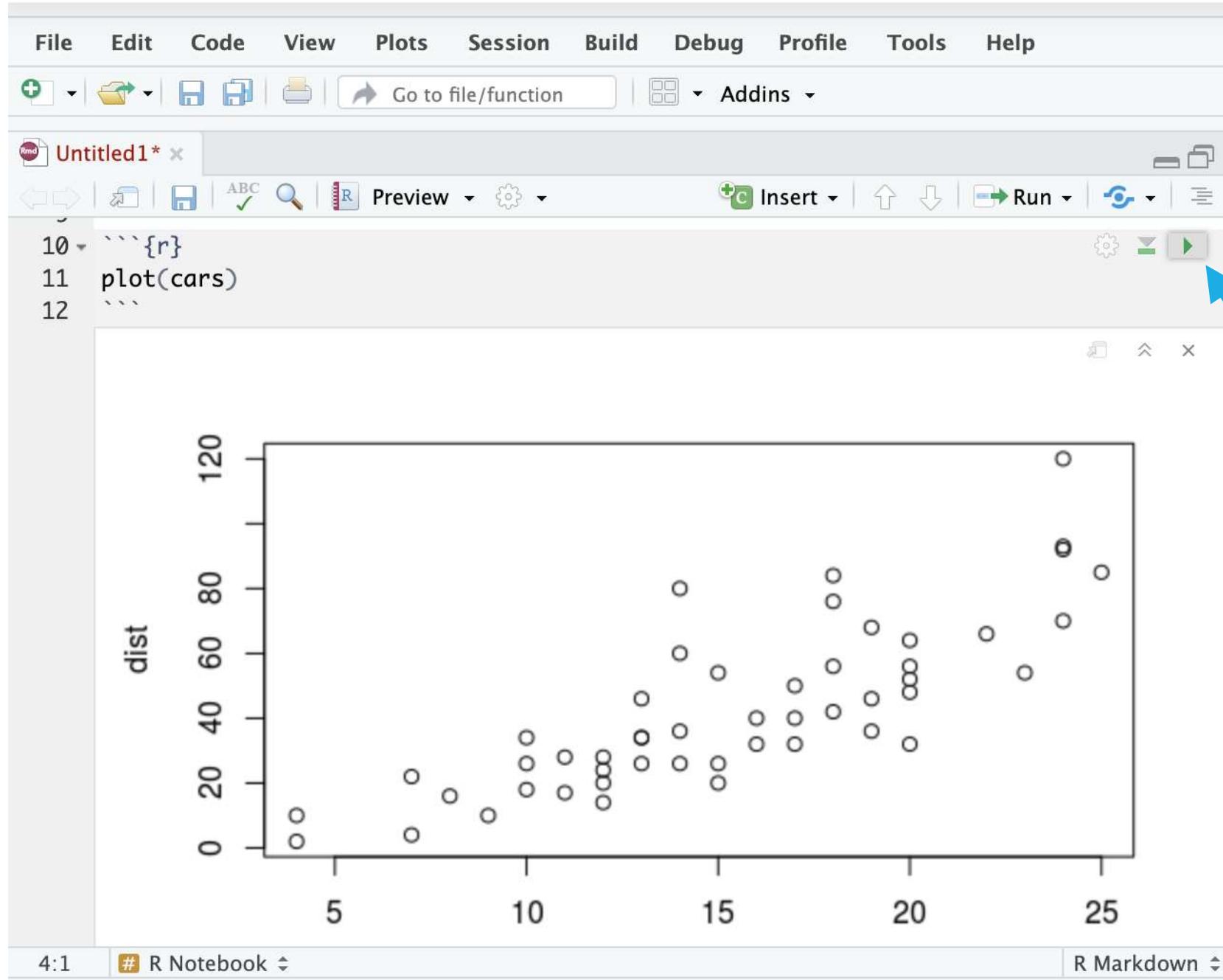
Chunk
name

Run chunk

Code in body
of chunk

```
18  ``{r cars}  
19  summary(cars)  
20  ...  
21
```





Run
code
chunk

Your Turn #2

Open Rstudio.Cloud

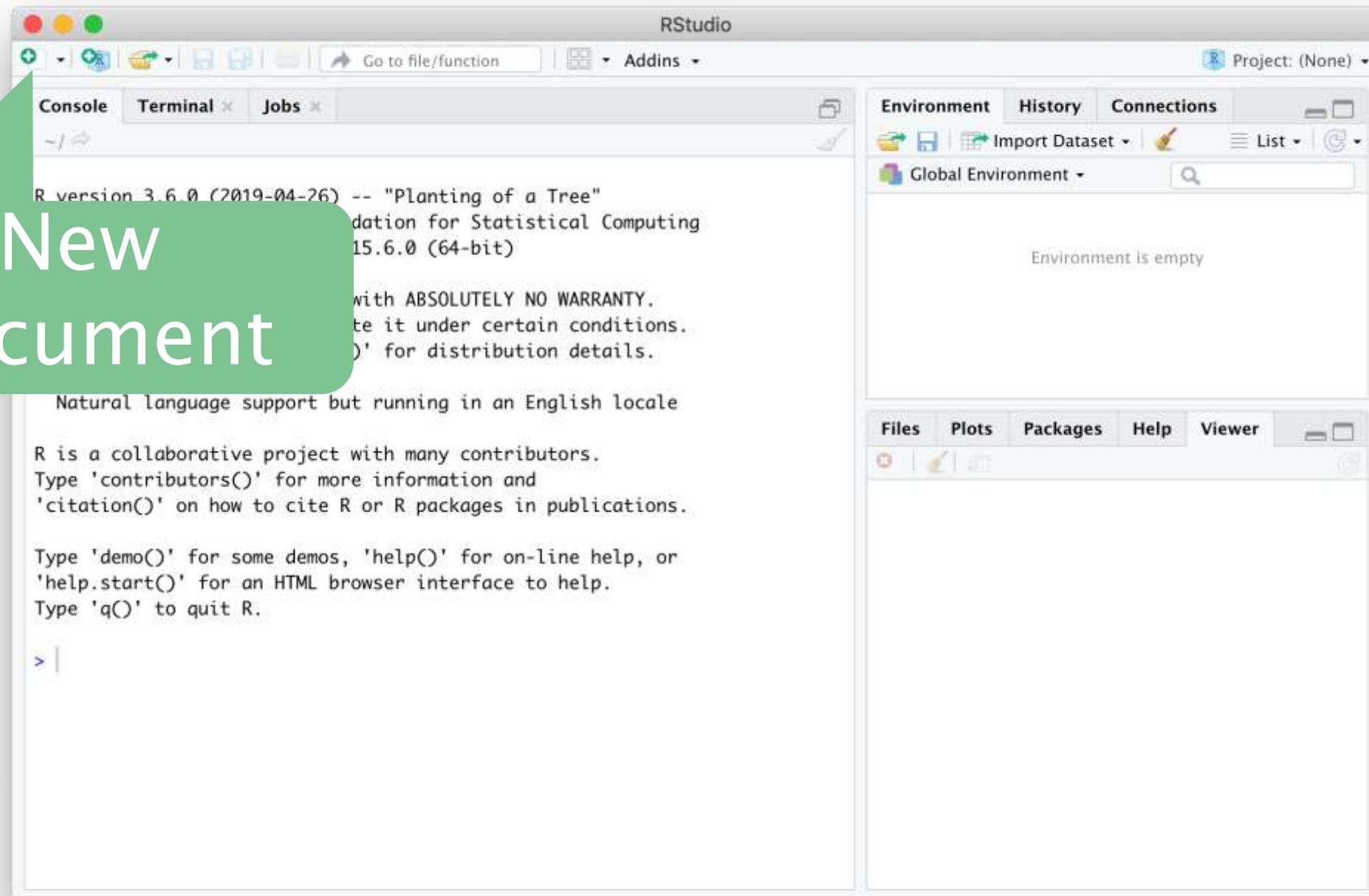
In the lower right MISC pane

Open 02 - Reproducible Reporting.

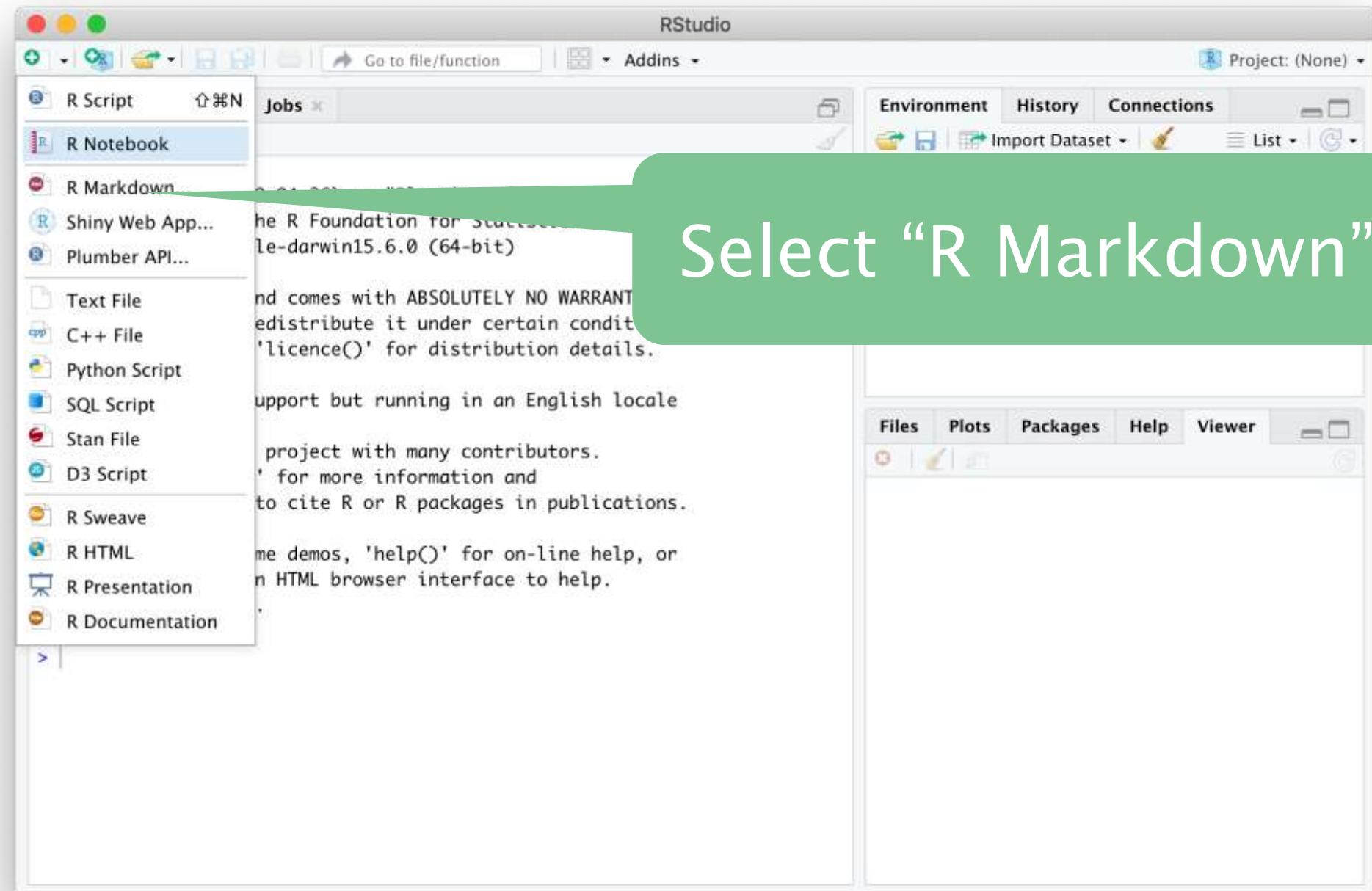
Read through the R Notebook and do everything it asks you to do.



Creating a New R Notebook



New Document



Select “R Markdown”

Save Document

```
1 ---  
2 title "Untitled"  
3 ---  
4  
5 ````{r}  
6 knitr::opts_chunk$set(echo = TRUE)  
7  
8
```

own document. Markdown is a simple formatting syntax for web pages, emails, and MS Word documents. For more details on using R Markdown, see <http://rmarkdown.rstudio.com>.

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:

```
14  
15 ````{r cars}  
16 summary(cars)  
17  
18  
19 ## Including Plots  
20  
21 You can also embed plots, for example:  
22  
23 ````{r pressure, echo=FALSE}  
24 plot(pressure)  
25
```

2:18 # Untitled R Markdown

header

code chunk

Untitled

R Markdown

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>.

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:

```
summary(cars)
```

```
##      speed      dist
## Min.   :4.0   Min.   : 2.00
## 1st Qu.:12.0  1st Qu.: 26.00
## Median :15.0  Median : 36.00
## Mean   :15.4  Mean   : 42.98
## 3rd Qu.:19.0  3rd Qu.: 56.00
## Max.   :25.0  Max.   :120.00
```

Including Plots

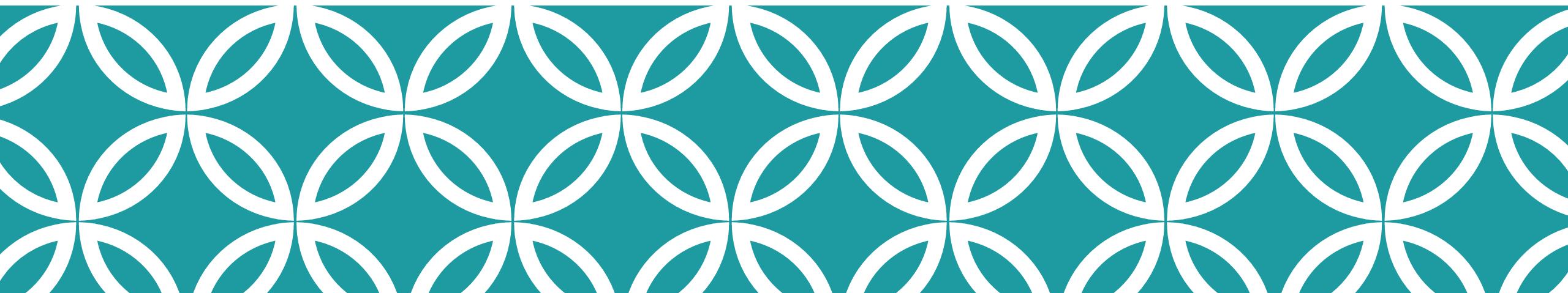
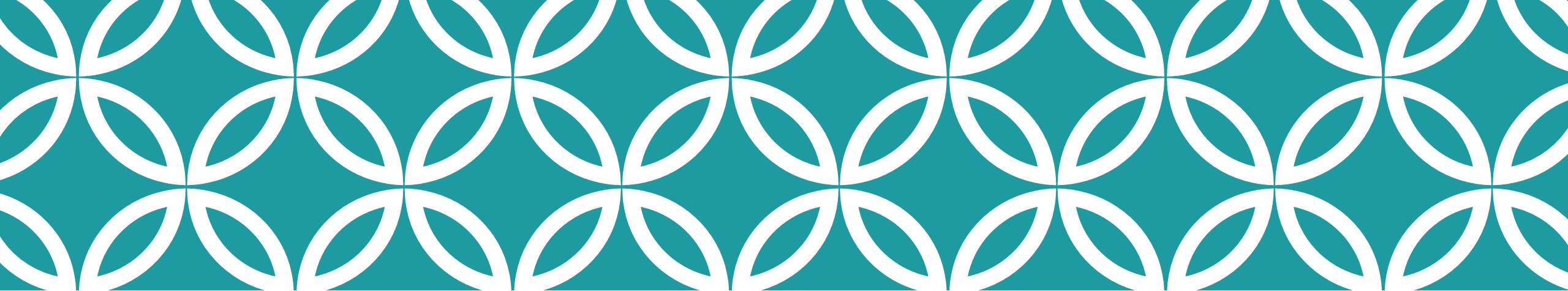
You can also embed plots, for example:

Your Turn #2

Try to reproduce this markdown document starting with the file you just created as a template.

The screenshot shows a browser window displaying an R Markdown document. The title is "My Very First Rmarkdown Document" by Amrom. It includes a section titled "Why R Markdown?", which discusses the seamless comingling of text and code. A code chunk is shown with the output "## [1] 15". Below this, there's a section on "Flexible Formatting" with examples of quotations and links.





What Else?

Why name Code Chunks?

Jump
between
chunks

The screenshot shows an RStudio interface with a code editor and a console. In the code editor, there are several code chunks:

```
16 When you click the **Knit** button a document includes both content as well as the output of within the document. You can embed an R code c
17
18 - ``{r cars}
19 summary(cars)
20 ``
21
22 - ## Including Plots
23
24 Yo Sample Markdown s, for example:
25 Chunk 1: setup
26 - ``
27 R Markdown
28 pl Chunk 2: cars
29 ``
30 Including Plots
31 No Chunk 3: pressure
32 ``
33 # Sample Markdown
```

A tooltip is displayed over the word "Sample" in the line "Yo Sample Markdown". The tooltip contains the following text:

Sample Markdown
Chunk 1: setup
R Markdown
Chunk 2: cars
Including Plots
No Chunk 3: pressure

The console at the bottom shows the path "/cloud/project/" and the R version information "R version 3.5.2 (2018-12-20) -- \"Froshell Taloa\"".

“Setup” Chunk & Chunk Options

chunk name
(optional)

“chunk option”
don't show code in rendered document

```
6 - ```{r setup, include=FALSE}
7 library(tidyverse)
8 library(lubridate)
9 ...
```

for dealing
with dates

Chunk options

| option | default | effect |
|------------|----------|---|
| eval | TRUE | Whether to evaluate the code and include its results |
| echo | TRUE | Whether to display code along with its results |
| warning | TRUE | Whether to display warnings |
| error | FALSE | Whether to display errors |
| message | TRUE | Whether to display messages |
| tidy | FALSE | Whether to reformat code in a tidy way when displaying it |
| results | "markup" | "markup", "asis", "hold", or "hide" |
| cache | FALSE | Whether to cache results for future renders |
| comment | "##" | Comment character to preface results with |
| fig.width | 7 | Width in inches for plots created in chunk |
| fig.height | 7 | Height in inches for plots created in chunk |

Keyboard Shortcuts

Insert a code chunk into white space within your open R Markdown document using:

- Windows: CTRL+ALT+i
- Mac: COMMAND+OPTION+I

Execute using shortcuts:

- Windows: CTRL+SHIFT+ENTER
- Mac: COMMAND+SHIFT+ENTER

Multiple Output Formats are Available

Pandoc universal document converter can create multiple document types:

- html
- pdf
- docx

Can also create presentations and dashboards

- Including Powerpoint (most recent version of RStudio)

Creating a pdf report

R & RStudio require additional packages to create a nicely formatted pdf report

Behind the scenes, R will use a markup language called LaTeX to turn your markdown into the pdf

Install the tinytex package with the following commands:

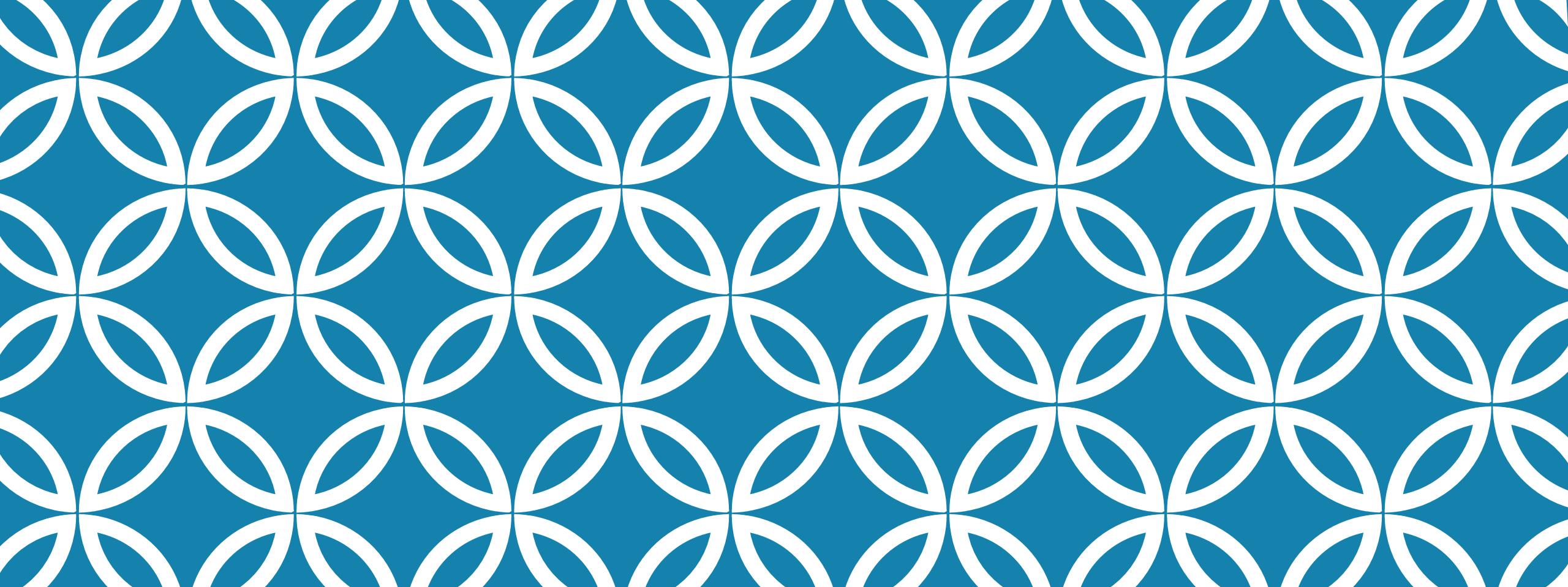
- `install.packages('tinytex')`
- `tinytex::install_tinytex()`

Goals

1. Understand why reproducible reporting is important
2. Learn to work within R Markdown for reproducible reports

Objectives

1. Create an R Markdown document and generate different types of output files
2. Practice modifying each component of a R Markdown file



Data Visualization

Session 3
Stephan Kadauke
July 16, 2020

| July 16 2020 | Session | Instructor |
|-------------------|--|-------------------|
| 1:00 pm - 1:30 pm | Instructor Introductions, Introduction to technology | Amrom Obstfeld |
| 1:30 pm - 2:15 pm | Introduction to R and RStudio | Joe Rudolf |
| 2:30 pm - 3:15 pm | Reproducible Reporting | Patrick Mathias |
| 3:30 pm - 5:00 pm | Data Visualization | Stephan Kadauke |
| July 17 2020 | | |
| 1:00 pm - 2:30 pm | Data Transformation | Amrom Obstfeld |
| 2:45 pm - 4:15 pm | Statistical Analysis | Dan Herman |
| 4:30 pm - 5:00 pm | Advanced Reporting | Patrick Mathias |

Goals

1. Appreciate the importance of visualization for understanding data
2. Learn how to use ggplot2 to visualize data

Objectives

1. Create a basic visualization using a simple template
2. Define “aesthetic mapping” and explain how aesthetic mappings relate variables of a data frame to features of graphic markings on a plot
3. Write the code to specify a type of plot and fine tune its appearance using “geom” functions
4. Explain how to add layers to a ggplot object to create complex and highly customized visualizations

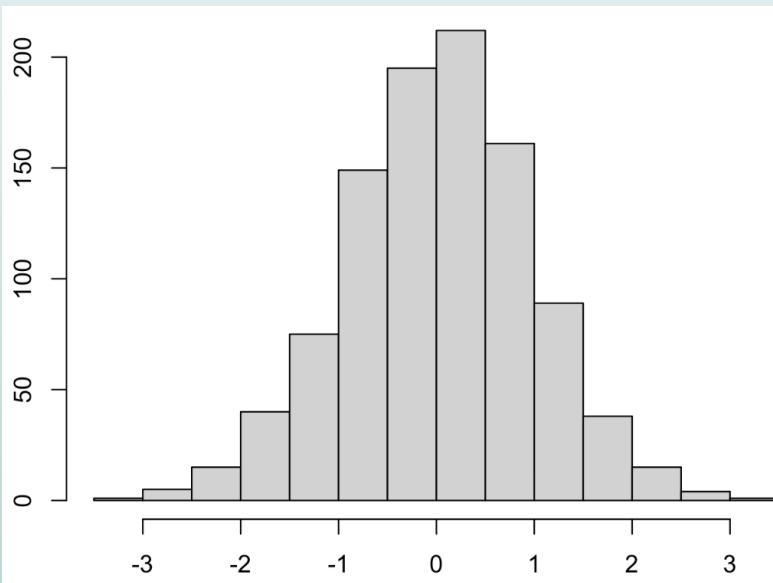
Your Turn 1

Consider the `covid_testing` data frame.

What do you think plot would look like in which:

- the x-axis represents `pan_day` (day of the pandemic), and
- the y-axis represents the number of tests that were performed on that day?

Your Turn 2



What is the name of this kind of plot?
Type the answer into the chat!

Your Turn 3

Type the following code in the RStudio console to make a graph.

Pay attention to the spelling, capitalization, and parentheses!

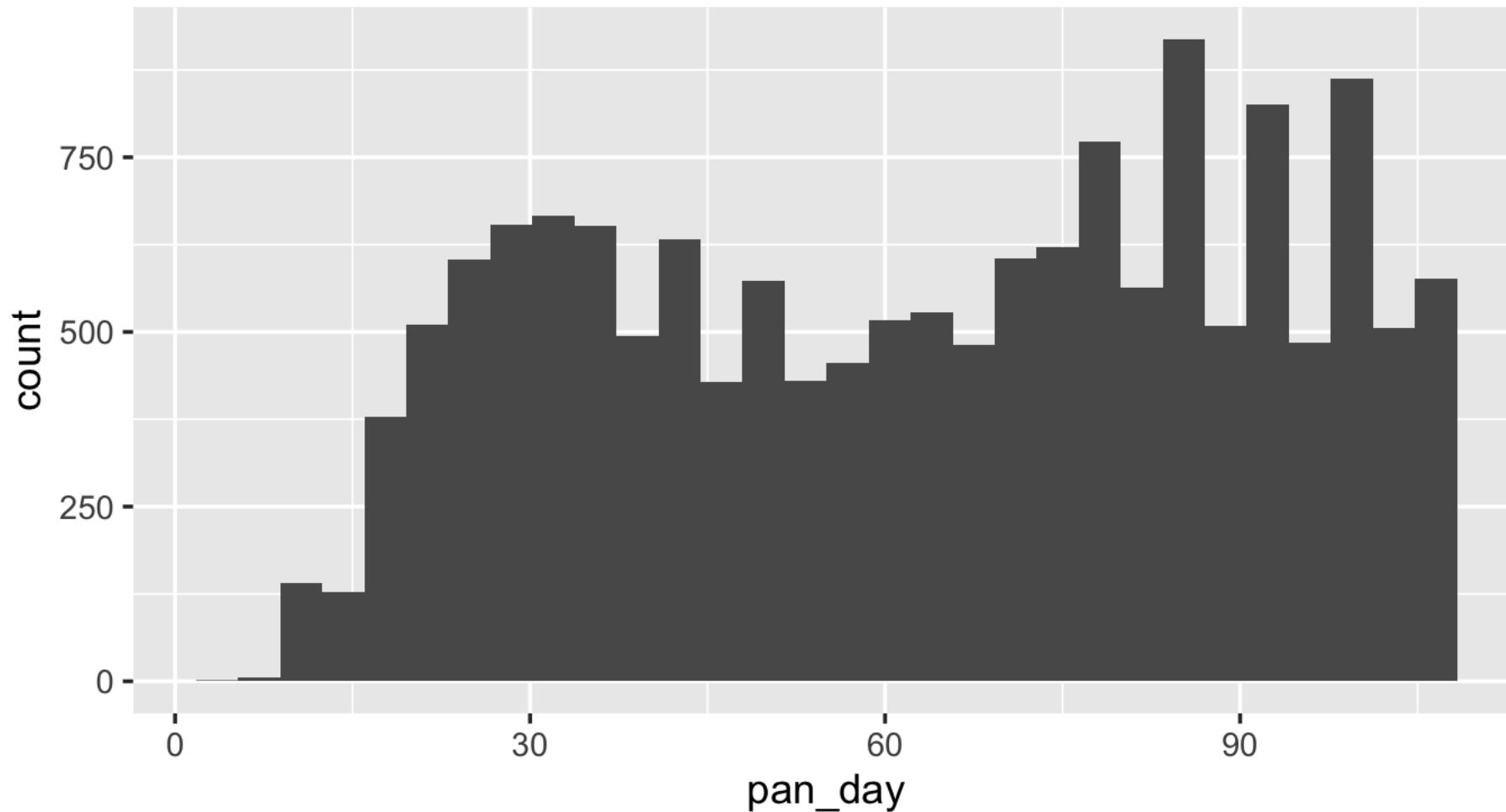
```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day))
```

When you run this code, you will get what looks like an error but is actually just a message.

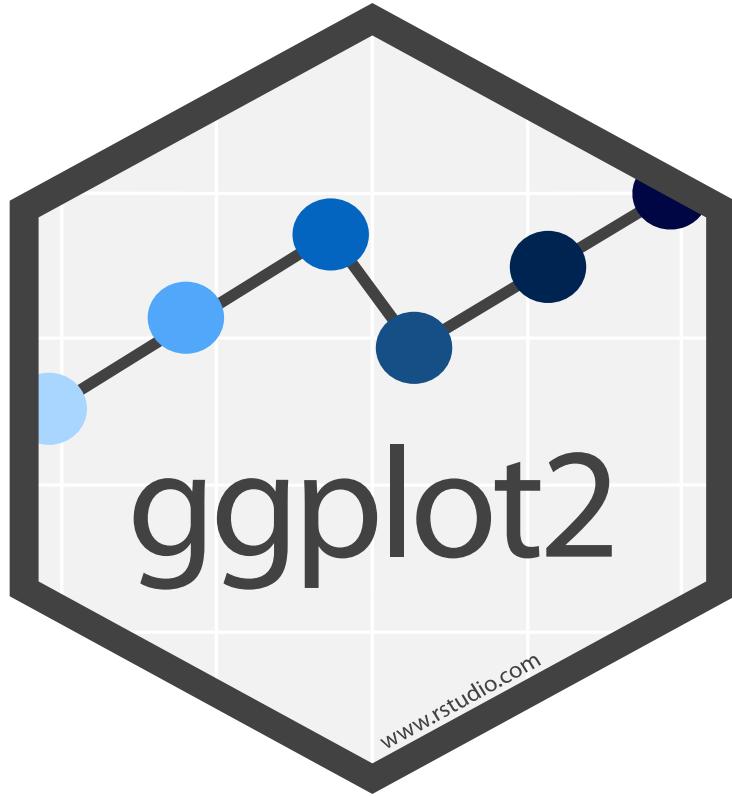
R lets you know that when you ask it to draw a histogram you should tell it how wide each bin should be, because this affects the granularity of the data displayed.

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

```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day))
```



```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day))
```



ggplot()

Always start
with ggplot()

data frame

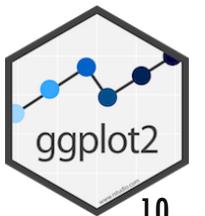
+ sign
before new line

```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day))
```

type of plot

mappings inside
aes() function

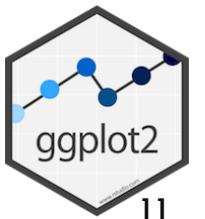
x axis
mapping



To make **any** kind of graph:

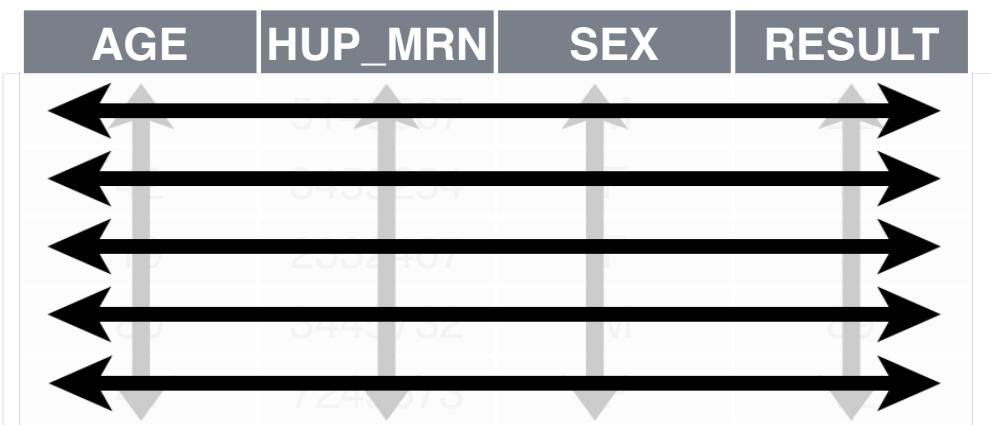
1. Choose a “tidy”
data frame

```
ggplot(data = data_frame) +  
  geom_function(mapping = aes(mappings))
```



1. Pick a “Tidy” Data Frame

| AGE | HUP_MRN | SEX | RESULT |
|-----|---------|-----|--------|
| 1 | 1 | 1 | 1 |
| 2 | 2 | 2 | 2 |
| 3 | 3 | 3 | 3 |
| 4 | 4 | 4 | 4 |



A data set is **tidy** if:

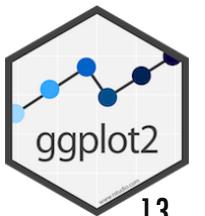
1. Each **variable** is in its own **column**
2. Each **observation** is in its own **row**
3. Each **value** is in its own **cell**

To make **any** kind of graph:

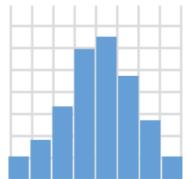
1. Choose a “tidy”
data frame

```
ggplot(data = data_frame) +  
  geom_function(mapping = aes(mappings))
```

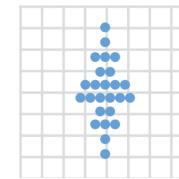
2. Pick a “**geom**”
function



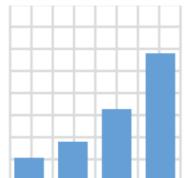
2. Choose a “Geom” Function



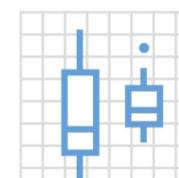
`geom_histogram()`



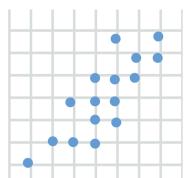
`geom_dotplot()`



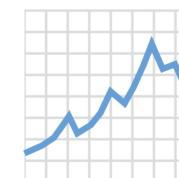
`geom_bar()`



`geom_boxplot()`



`geom_point()`



`geom_line()`



Data Visualization with ggplot2 :: CHEAT SHEET

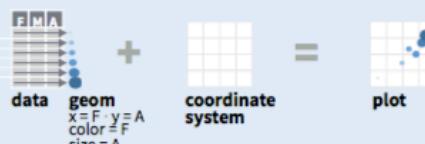


Basics

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



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.



Complete the template below to build a graph.

```
ggplot(data = <DATA>) +  
  <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>),  
    stat = <STAT>, position = <POSITION>) +  
  <COORDINATE_FUNCTION> +  
  <FACET_FUNCTION> +  
  <SCALE_FUNCTION> +  
  <THEME_FUNCTION>
```

Annotations for the template:

- required**: **GEOM_FUNCTION**
- Not required, sensible defaults supplied**: **COORDINATE_FUNCTION**, **FACET_FUNCTION**, **SCALE_FUNCTION**, **THEME_FUNCTION**

ggplot(data = mpg, aes(x = cty, y = hwy)) Begins a plot that you finish by adding layers to. Add one geom function per layer.

```
aesthetic mappings  data  geom  
ggplot(x = cty, y = hwy, data = mpg, geom = "point")
```

Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemploy))  
b <- ggplot(seals, aes(x = long, y = lat))
```

a + geom_blank()
(Useful for expanding limits)

**b + geom_curve(aes(yend = lat + 1,
xend = long + 1, curvature = z)) - x, xend, y, yend,
alpha, angle, color, curvature, linetype, size**

**a + geom_path(lineend = "butt", linejoin = "round",
linemetre = 1)
x, y, alpha, color, group, linetype, size**

**a + geom_polygon(aes(group = group))
x, y, alpha, color, fill, group, linetype, size**

**b + geom_rect(aes(xmin = long, ymin = lat, xmax =
long + 1, ymax = lat + 1)) - xmax, xmin, ymax,
ymin, alpha, color, fill, linetype, size**

**a + geom_ribbon(aes(ymin = unemploy - 900,
ymax = unemploy + 900)) - x, ymax, ymin,
alpha, color, fill, group, linetype, size**

LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

**b + geom_abline(aes(intercept = 0, slope = 1))
b + geom_hline(aes(yintercept = lat))
b + geom_vline(aes(xintercept = long))**

**b + geom_segment(aes(yend = lat + 1, xend = long + 1))
b + geom_spoke(aes(angle = 1:1155, radius = 1))**

ONE VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
```

**c + geom_area(stat = "bin")
x, y, alpha, color, fill, linetype, size**

**c + geom_density(kernel = "gaussian")
x, y, alpha, color, fill, group, linetype, size, weight**

**c + geom_dotplot()
x, y, alpha, color, fill**

TWO VARIABLES

continuous x , continuous y

```
e <- ggplot(mpg, aes(cty, hwy))  
  
e + geom_label(aes(label = cty), nudge_x = 1,  
nudge_y = 1, check_overlap = TRUE) x, y, label,  
alpha, angle, color, family, fontface, hjust,  
lineheight, size, vjust
```

**e + geom_jitter(height = 2, width = 2)
x, y, alpha, color, fill, shape, size**

**e + geom_point(), x, y, alpha, color, fill, shape,
size, stroke**

**e + geom_quantile(), x, y, alpha, color, group,
linetype, size, weight**

**e + geom_rug(sides = "bl") x, y, alpha, color,
linetype, size**

**e + geom_smooth(method = lm), x, y, alpha,
color, fill, group, linetype, size, weight**

**e + geom_text(aes(label = cty), nudge_x = 1,
nudge_y = 1, check_overlap = TRUE) x, y, label,
alpha, angle, color, family, fontface, hjust,
lineheight, size, vjust**

discrete x , continuous y

```
f <- ggplot(mpg, aes(class, hwy))
```

**f + geom_col(), x, y, alpha, color, fill, group,
linetype, size**

**f + geom_boxplot(), x, y, lower, middle, upper,
ymax, ymin, alpha, color, fill, group, linetype,
shape, size, weight**

**f + geom_dotplot(binaxis = "y", stackdir =
"center"), x, y, alpha, color, fill, group**

**f + geom_violin(scale = "area"), x, y, alpha, color,
fill, group, linetype, size, weight**

continuous bivariate distribution

```
h <- ggplot(diamonds, aes(carat, price))  
  
h + geom_bin2d(binwidth = c(0.25, 500))  
x, y, alpha, color, fill, linetype, size, weight
```

**h + geom_density2d()
x, y, alpha, colour, group, linetype, size**

**h + geom_hex()
x, y, alpha, colour, fill, size**

continuous function

```
i <- ggplot(economics, aes(date, unemploy))
```

**i + geom_area()
x, y, alpha, color, fill, linetype, size**

**i + geom_line()
x, y, alpha, color, group, linetype, size**

**i + geom_step(direction = "hv")
x, y, alpha, color, group, linetype, size**

visualizing error

```
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)  
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))  
  
j + geom_crossbar(fatten = 2)  
x, y, ymax, ymin, alpha, color, fill, group, linetype,  
size
```

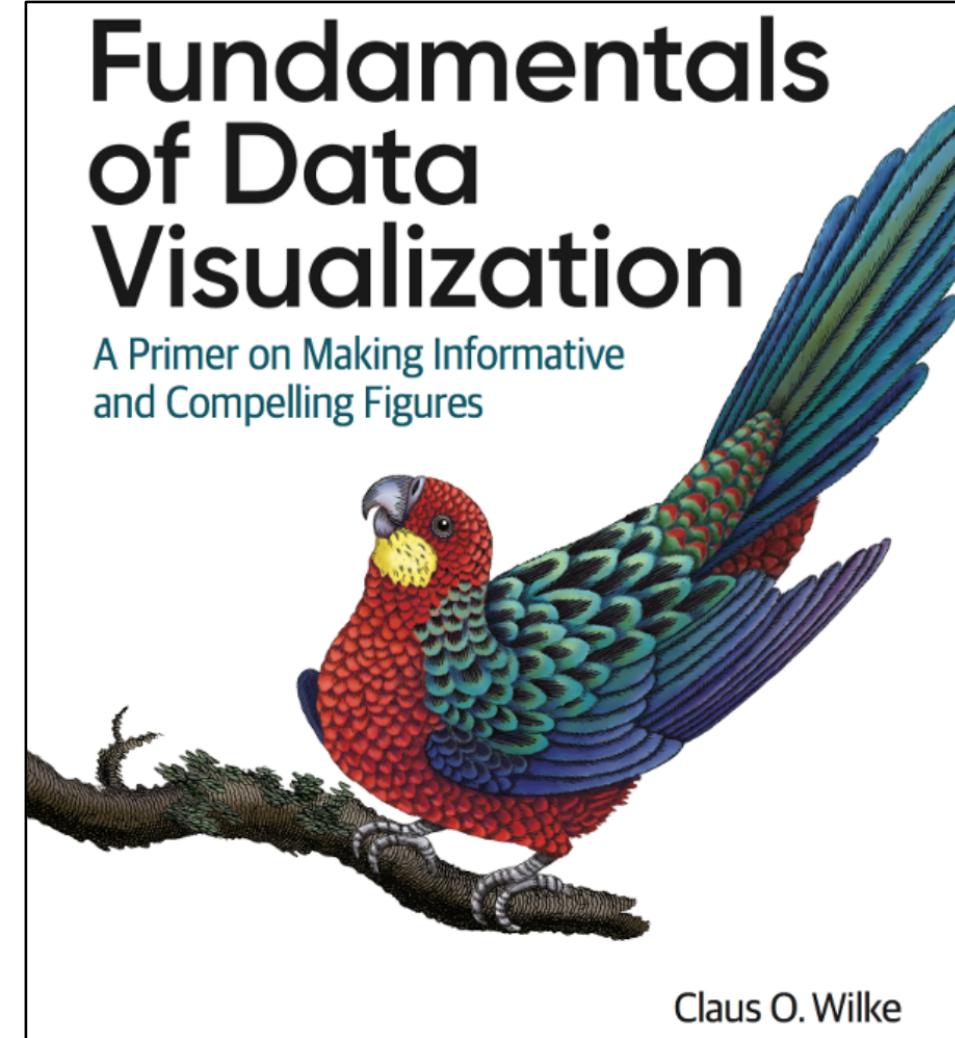
**j + geom_errorbar(), x, y, max, min, alpha, color,
group, linetype, size, width (also
geom_errorbarh())**

**j + geom_linerange()
x, y, min, max, alpha, color, group, linetype, size**

**j + geom_pointrange()
x, y, min, max, alpha, color, fill, group, linetype,
shape, size**

maps

```
data <- data.frame(murder = USArrests$Murder)
```



<https://serialmentor.com/dataviz>

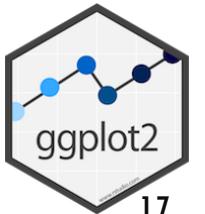
To make **any** kind of graph:

1. Choose a “tidy”
data frame

```
ggplot(data = data_frame) +  
  geom_function(mapping = aes(mappings))
```

2. Pick a “**geom**”
function

3. Write aesthetic
mappings



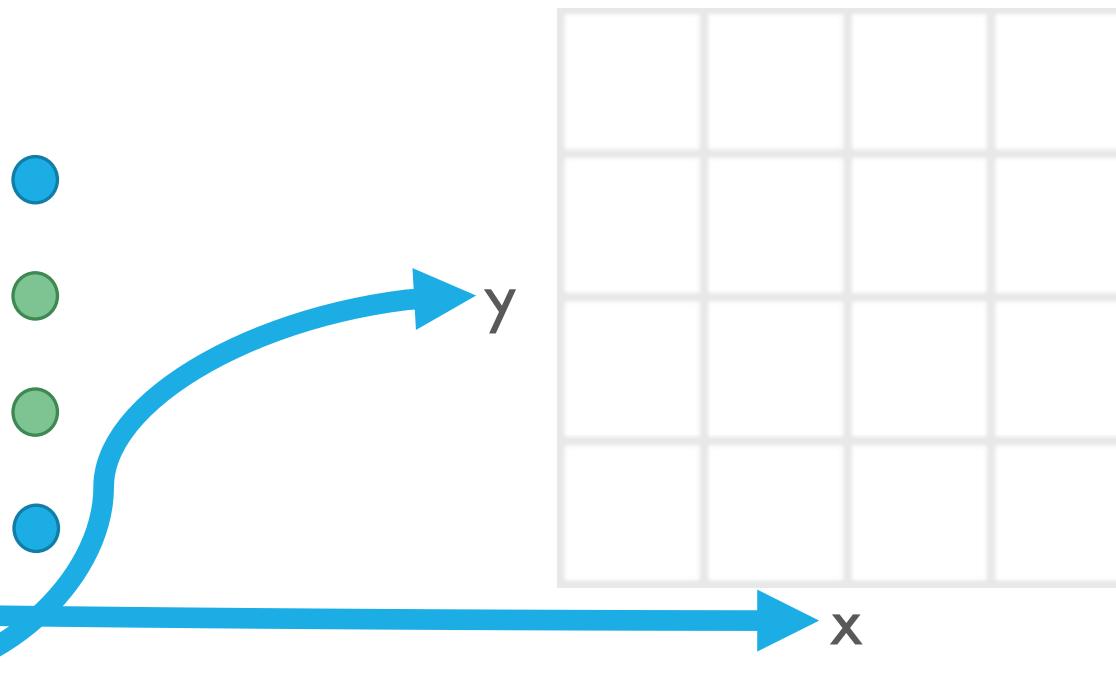
3. Write Aesthetic Mappings

```
aes(x = a, y = b, color = c)
```

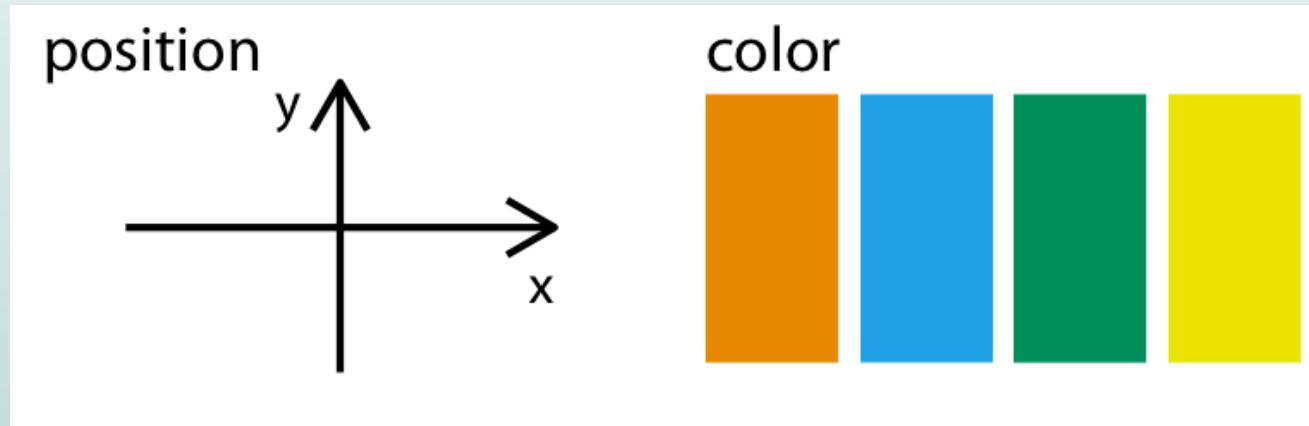
Data frame

| a | b | c |
|---|---|---|
| 1 | 3 | M |
| 2 | 1 | F |
| 3 | 3 | F |
| 2 | 2 | M |

Graph



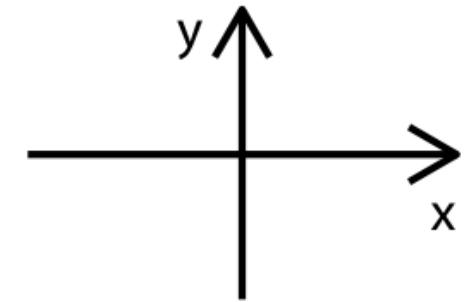
Your Turn 4



In addition to x/y position and color, what other aesthetics can you think of?
Type your answers in the chat!

Aesthetics

position



shape



size



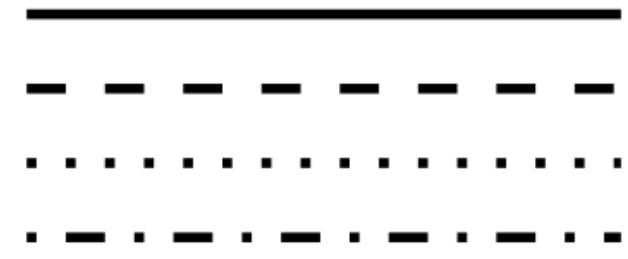
color



line width



line type



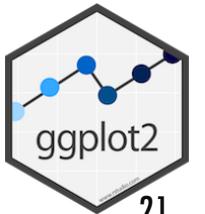
To make **any** kind of graph:

1. Choose a “tidy”
data frame

```
ggplot(data = data_frame) +  
  geom_function(mapping = aes(mappings))
```

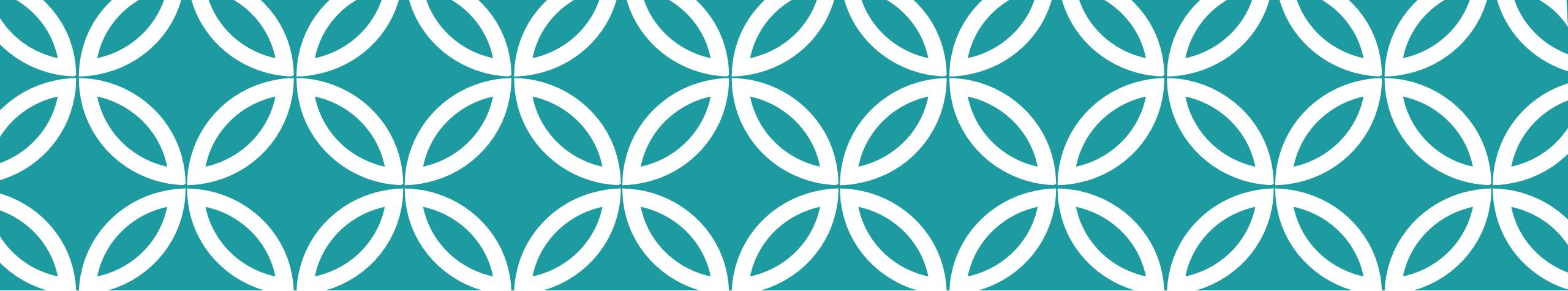
2. Pick a “**geom**”
function

3. Write aesthetic
mappings

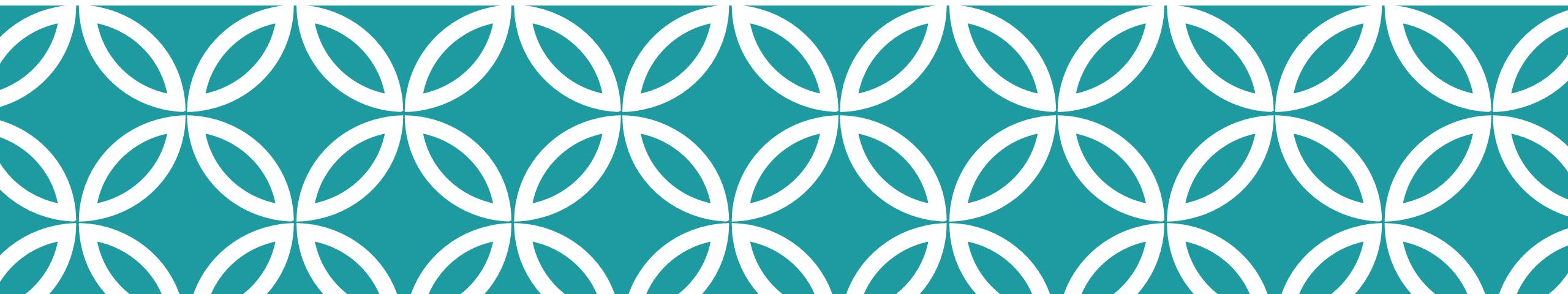


Your Turn 5

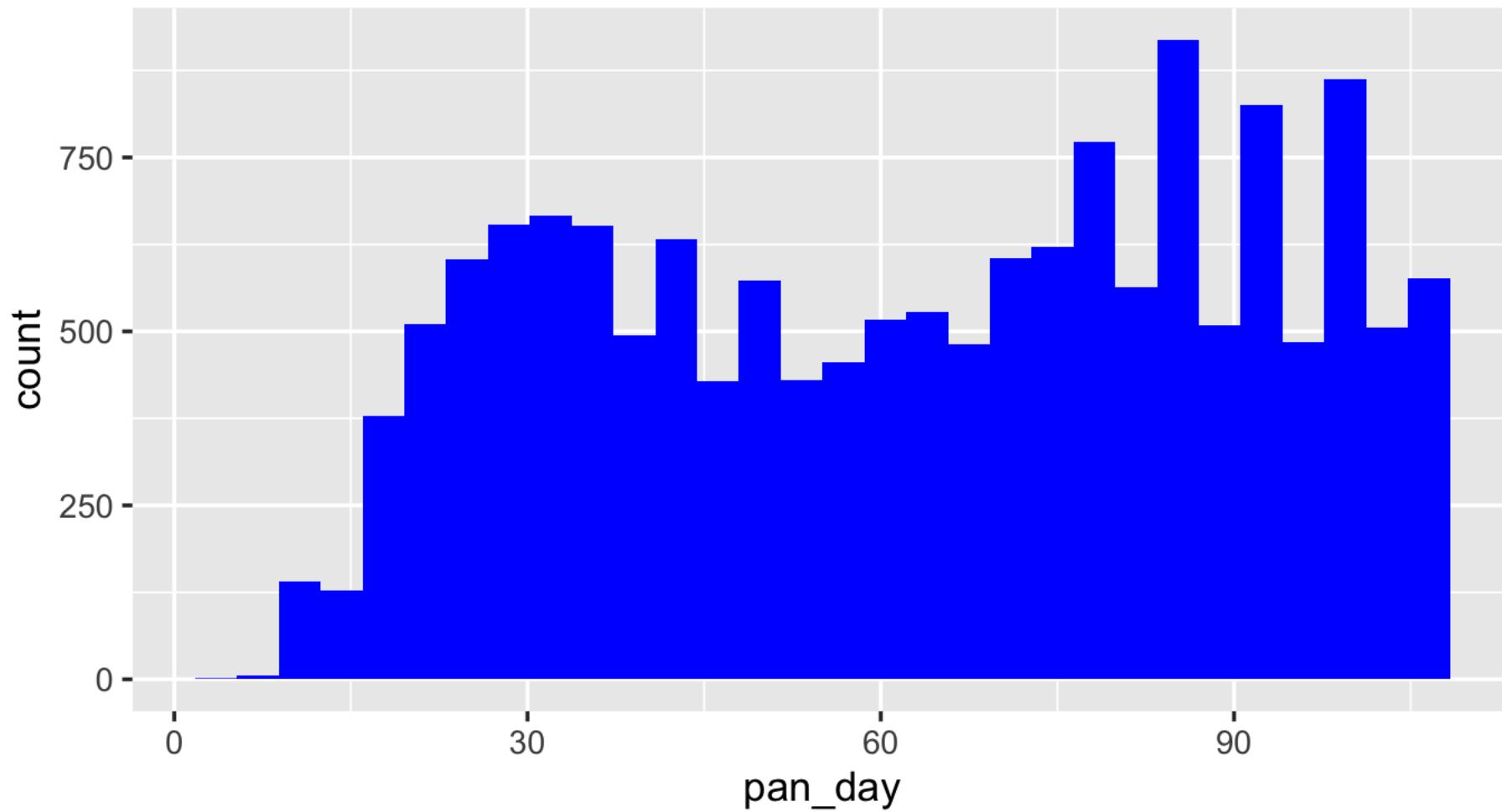
Open 03-Visualize.Rmd. Work through the exercises of the section titled “Your Turn 5.”

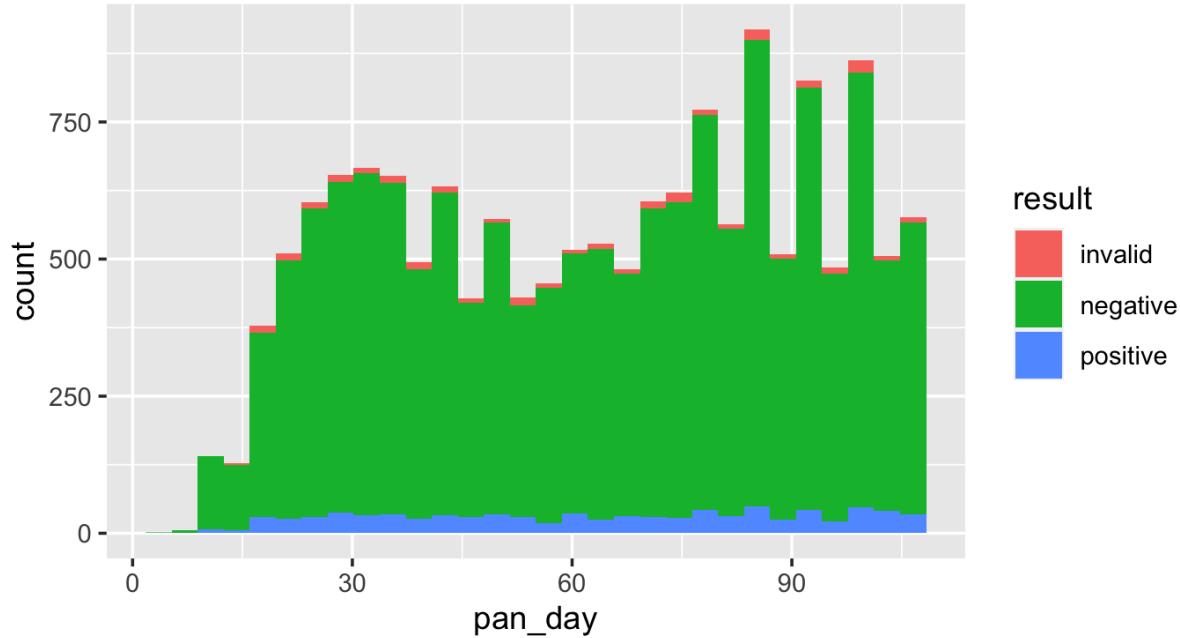


Setting vs **Mapping** Aesthetics



How would you make this plot?

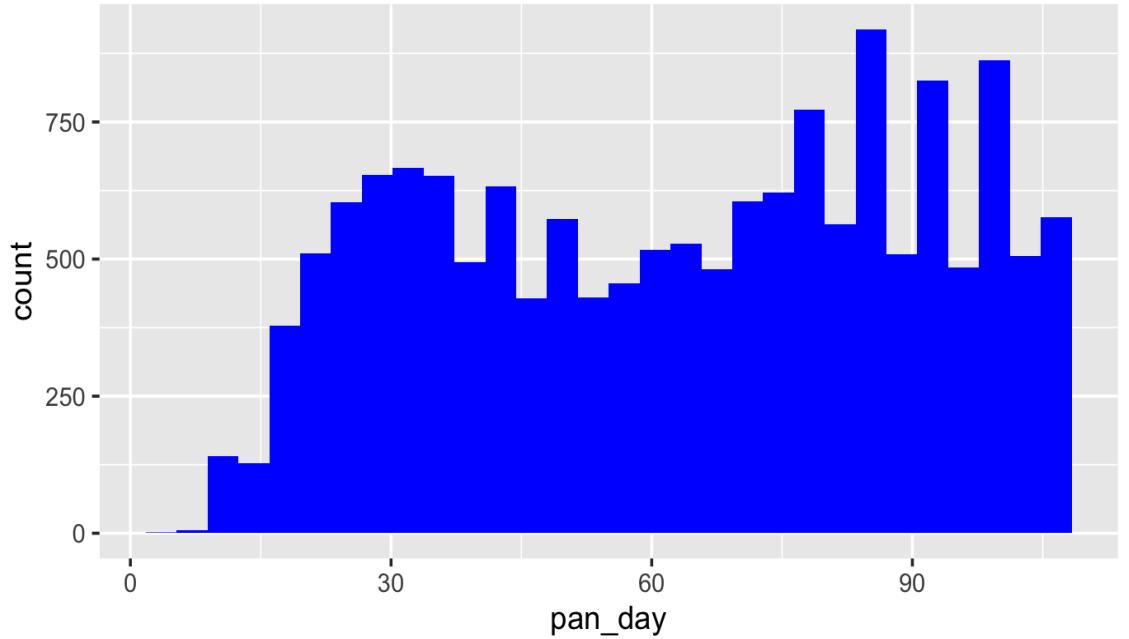




Inside of `aes()`:
map an aesthetic
to a variable

```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day, fill = result))
```

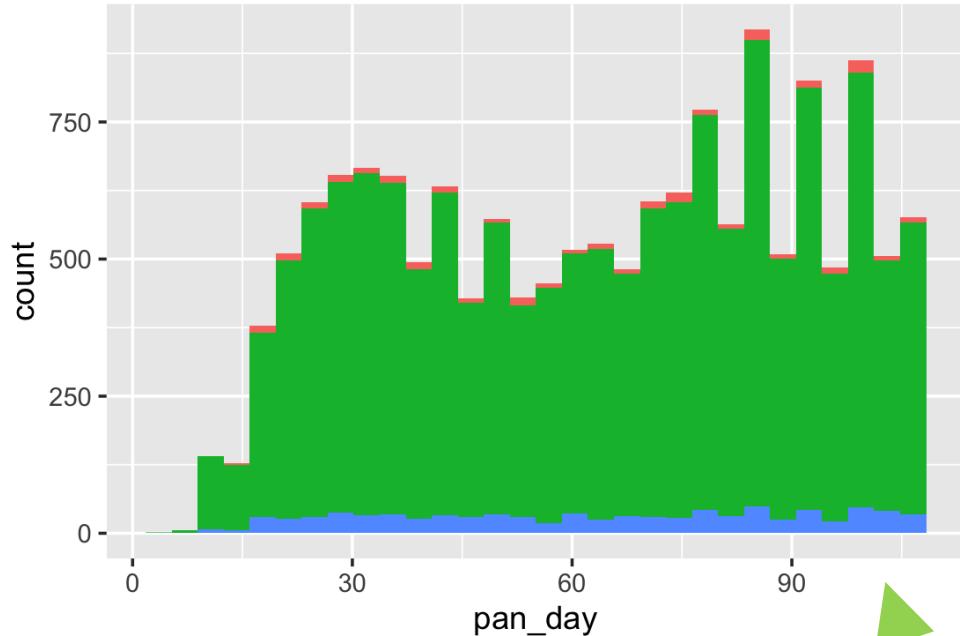




Outside of aes():
set an aesthetic to
a value

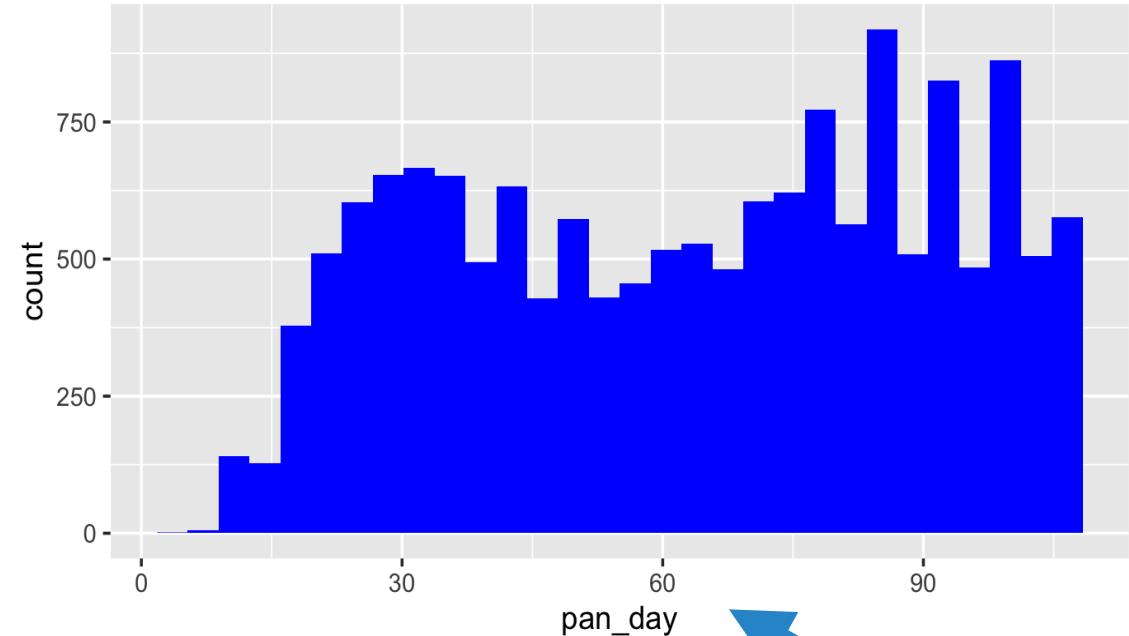
color name in
“quotes”

```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day), fill = "blue")
```



result

- invalid
- negative
- positive

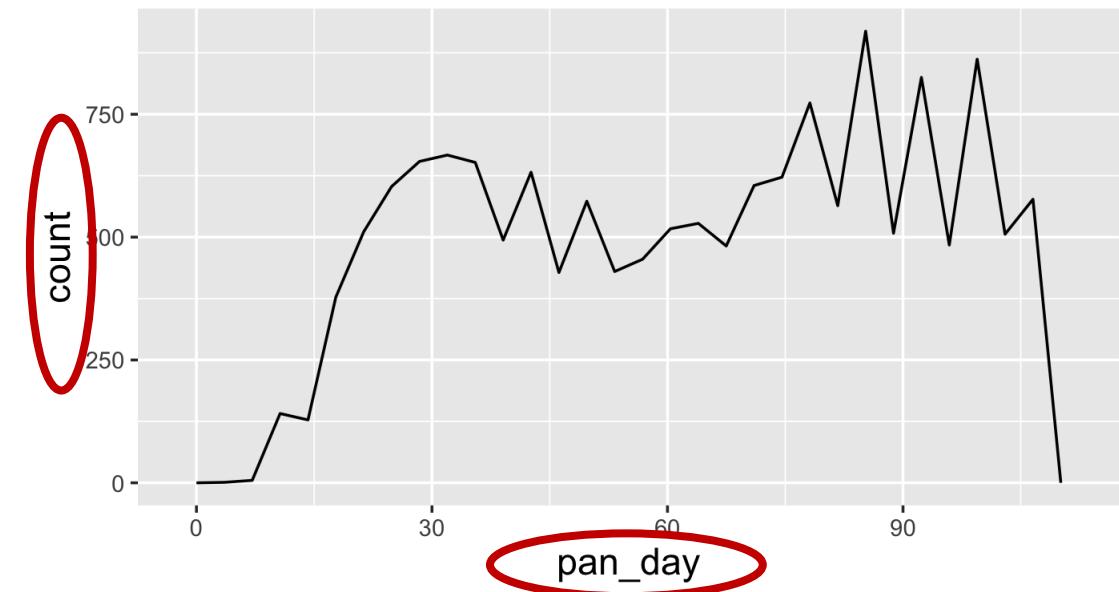
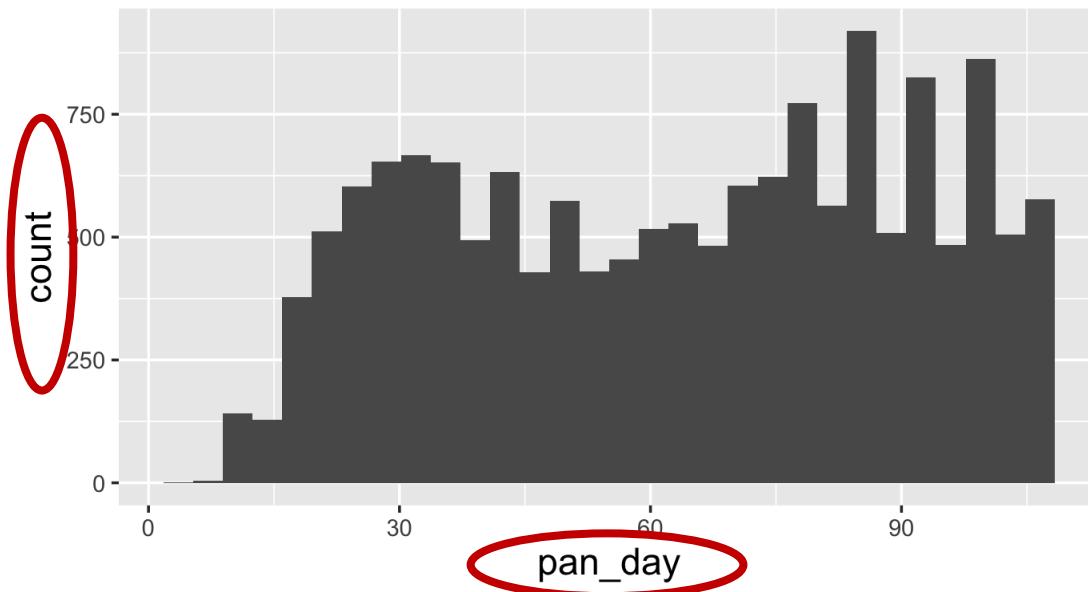


```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day, fill = result))
```

```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day), fill = "blue")
```

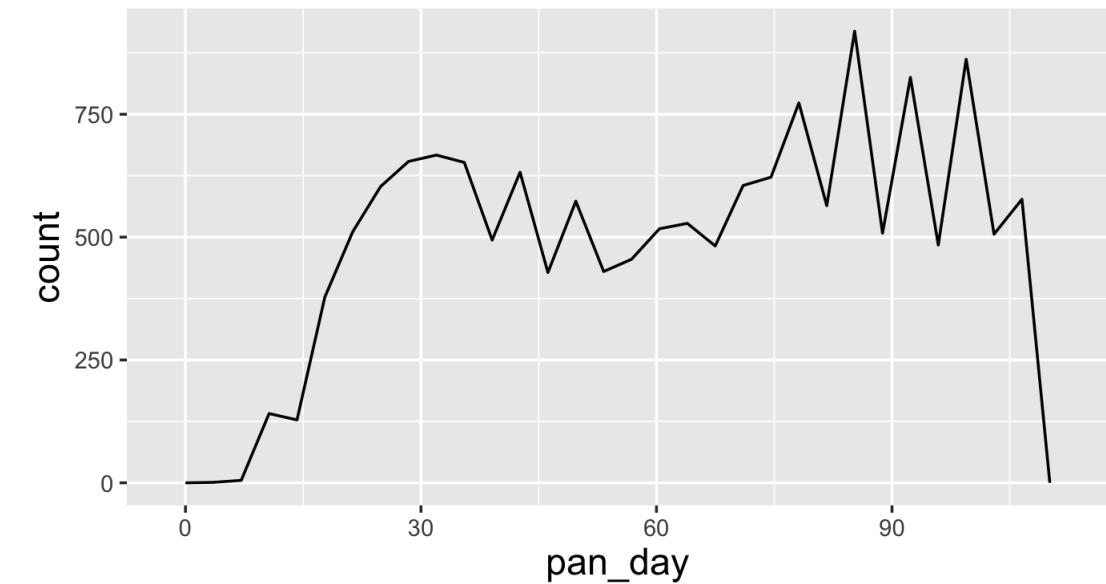
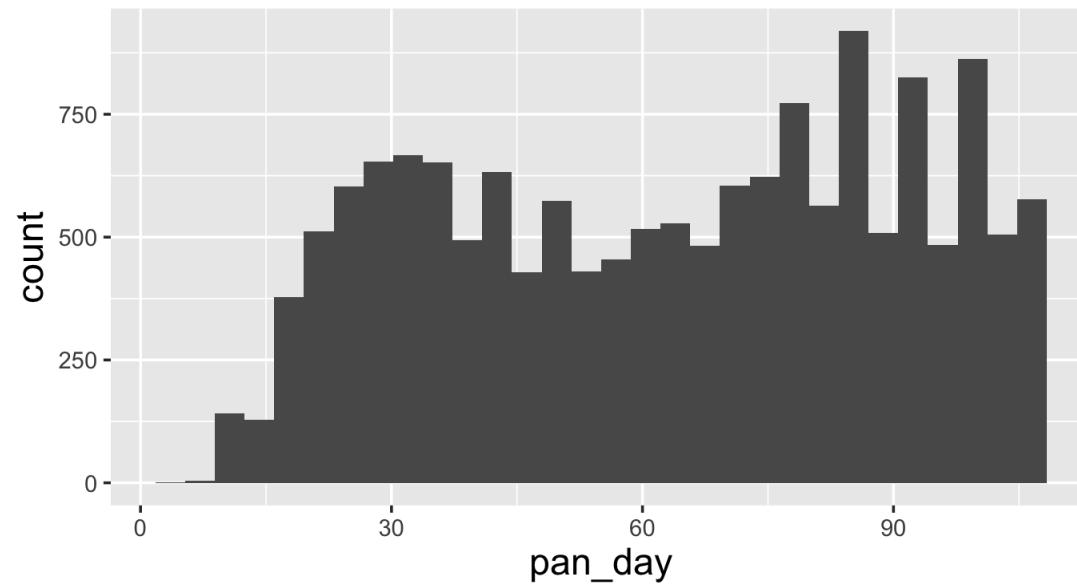
Geom Functions

How are these plots similar?



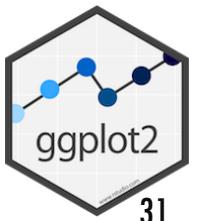
Same: x axis, y axis, data

How are these plots different?



Different **geometric object** (“geom”) used to represent the data

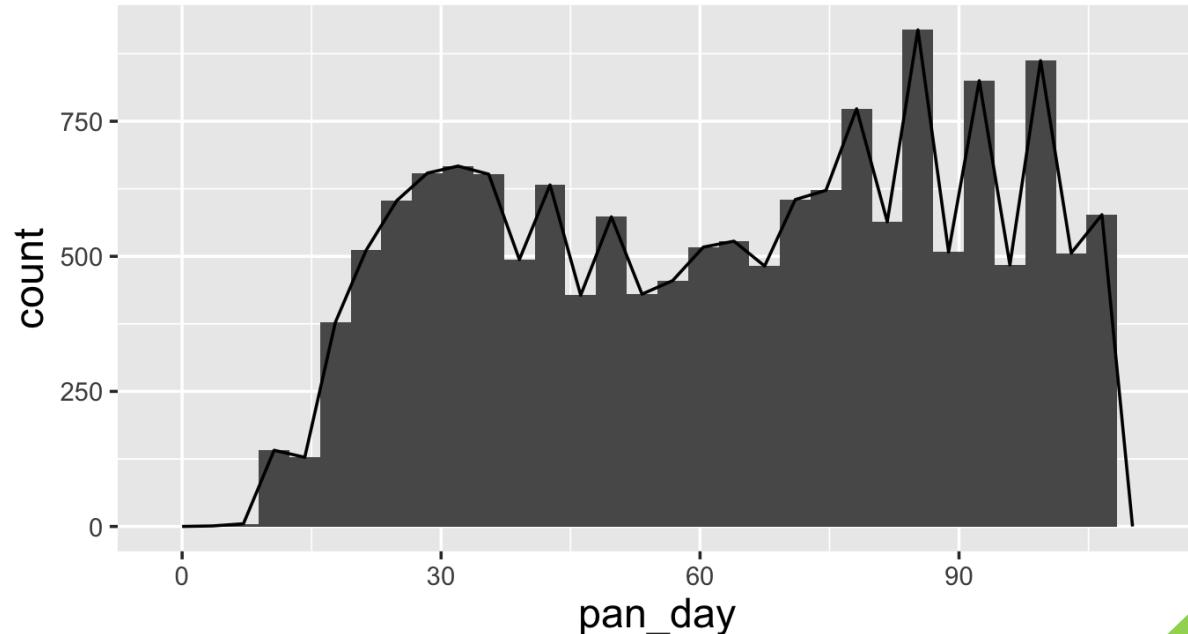
```
ggplot(data = data_frame) +  
  geom_function(mapping = aes(mappings))
```



Your Turn 6

Open 03-Visualize.Rmd. Work through the exercises of the section titled “Your Turn 6.”

Global vs Local Settings



Inside of `geom_` function:
apply **locally** to only the
current layer

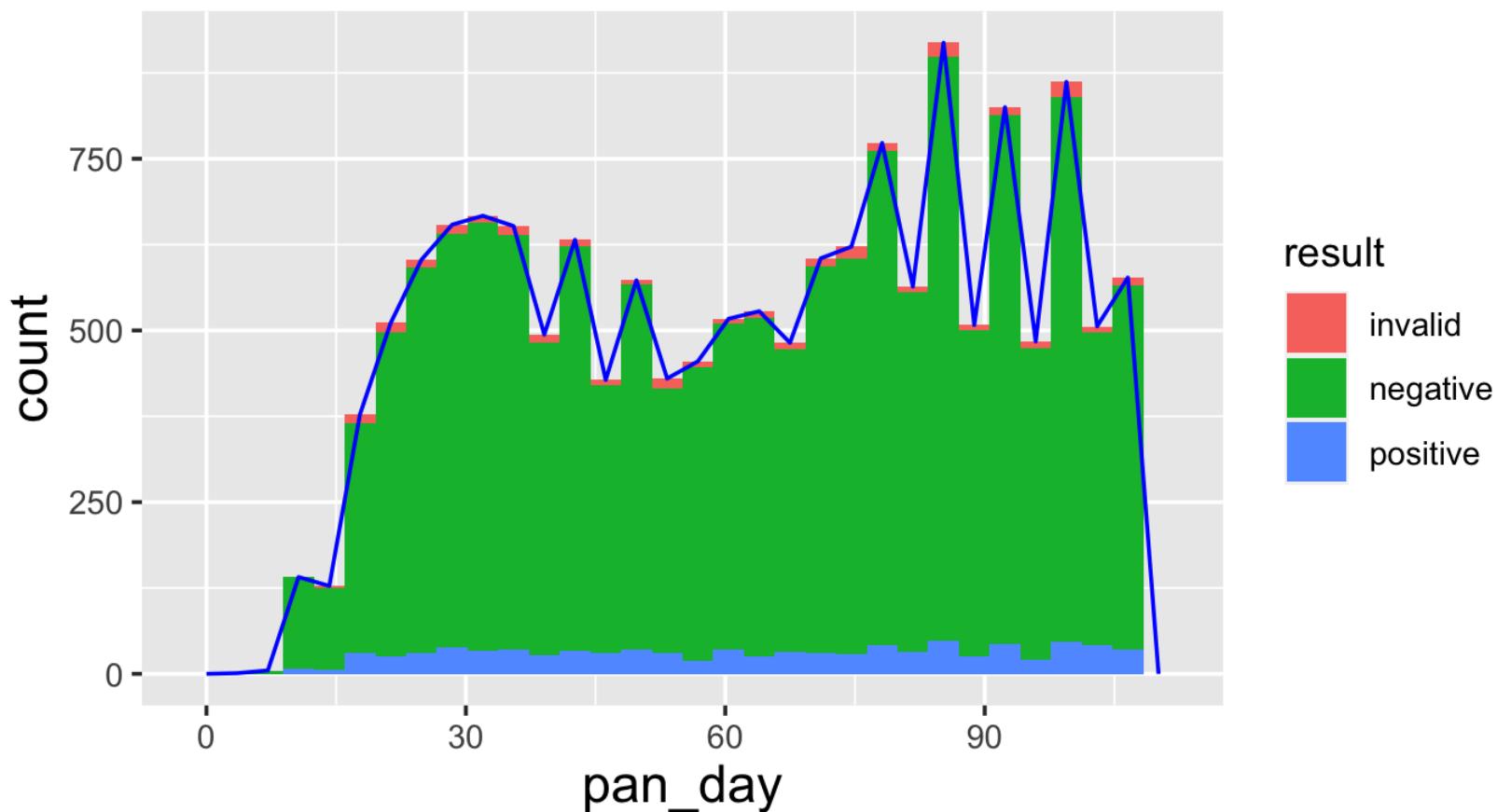
Inside of `ggplot()`:
apply **globally** to
every layer

```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day)) +  
  geom_freqpoly(mapping = aes(x = pan_day))
```

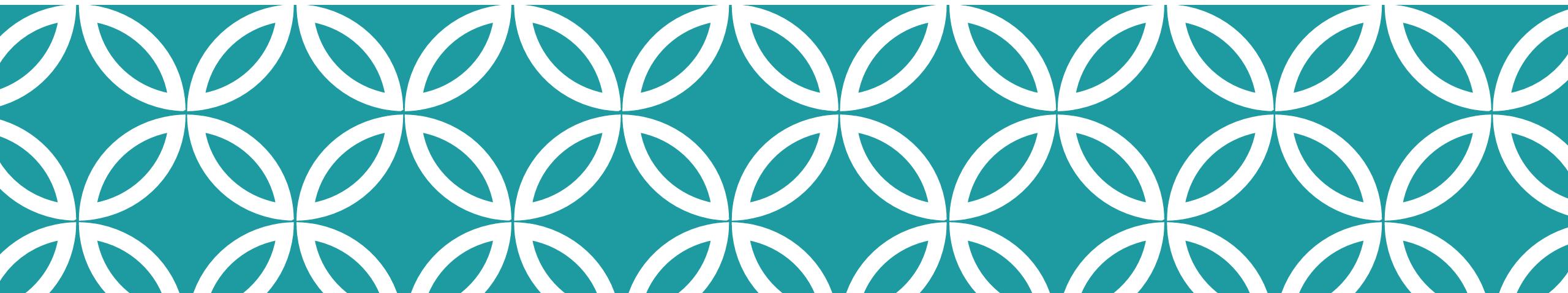
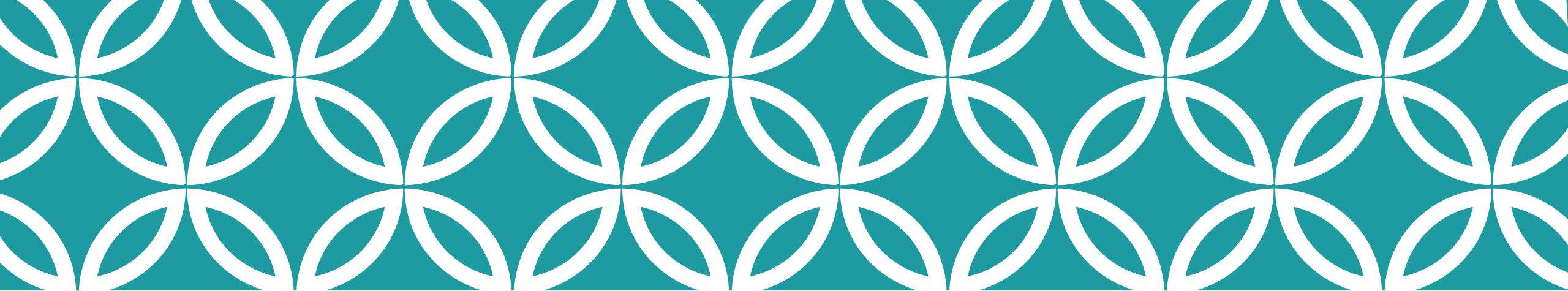
```
ggplot(data = covid_testing, mapping = aes(x = pan_day)) +  
  geom_histogram() +  
  geom_freqpoly()
```

empty ()

How would you make this plot?

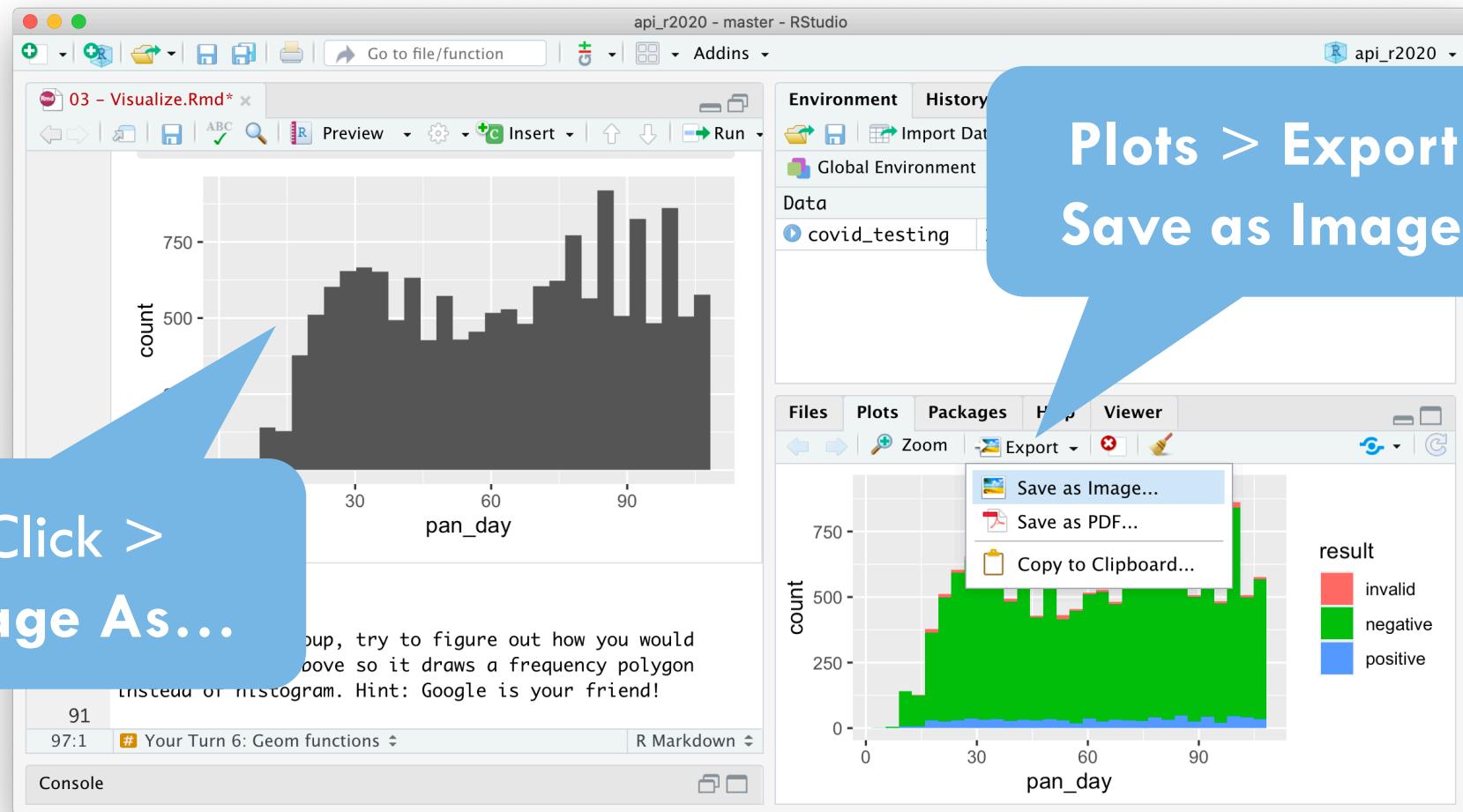


```
ggplot(data = covid_testing, mapping = aes(x = pan_day)) +  
  geom_histogram(mapping = aes(fill = result)) +  
  geom_freqpoly(color = "blue")
```



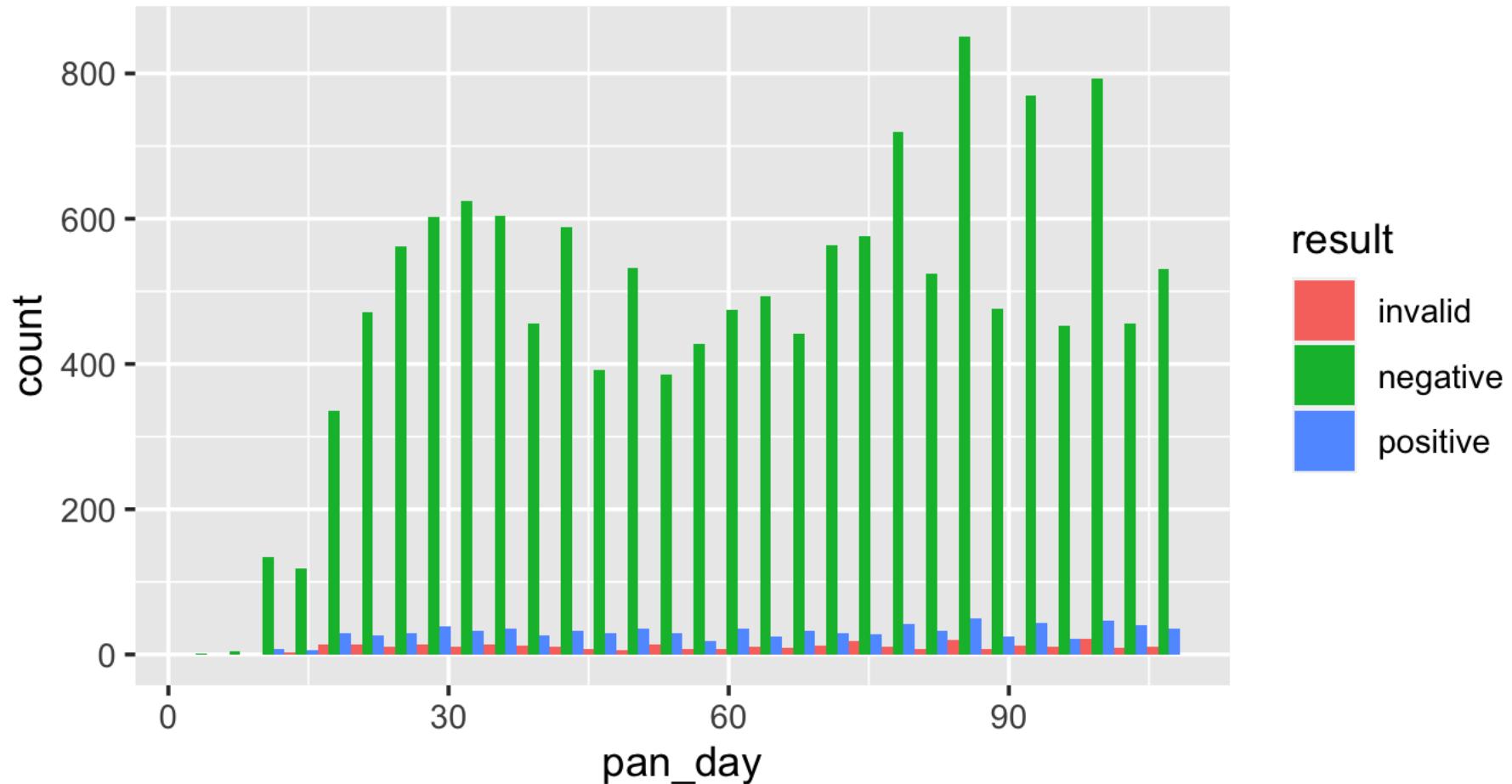
What Else?

Manually saving plots



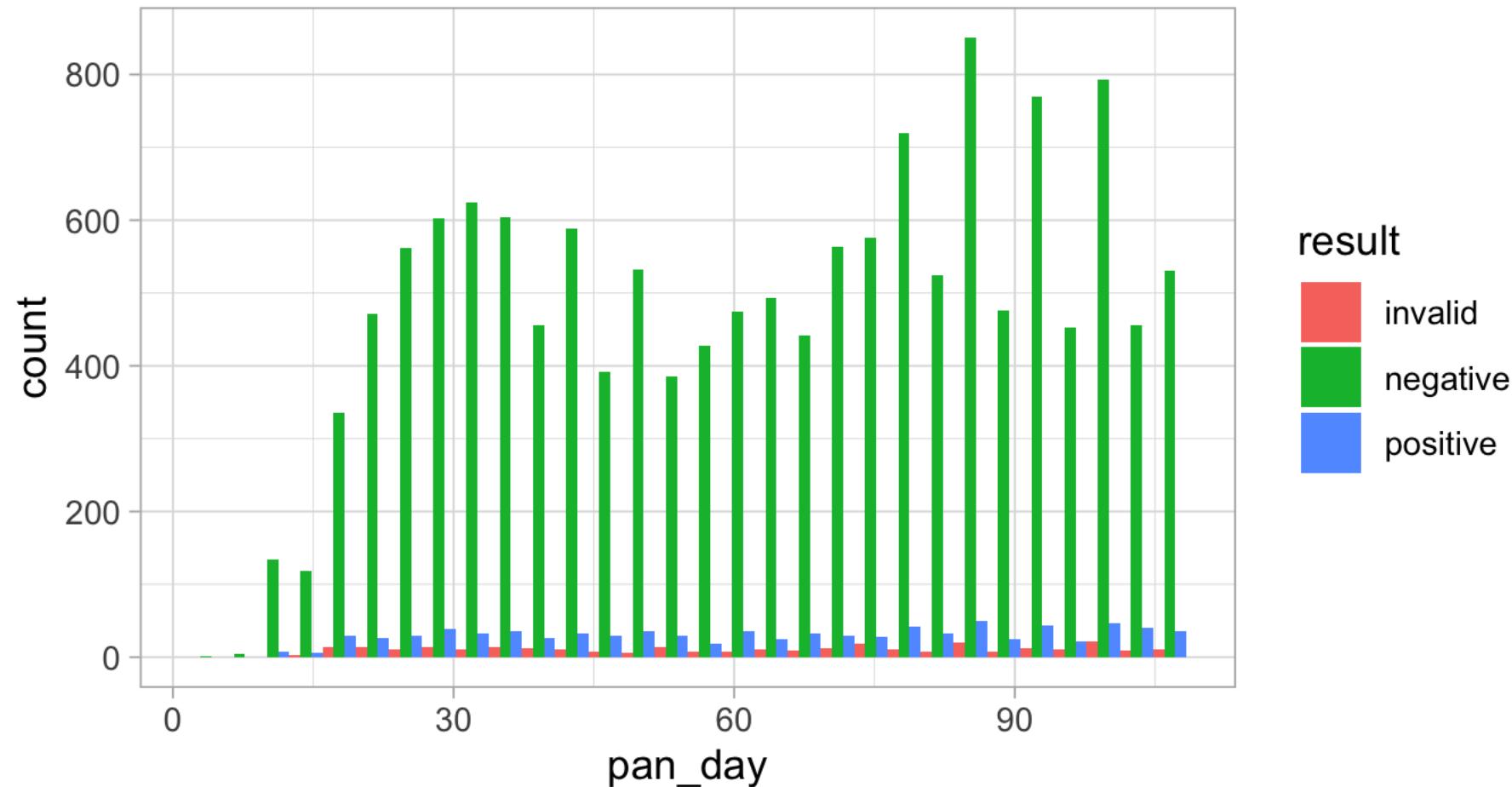
Position adjustments

How overlapping objects are arranged



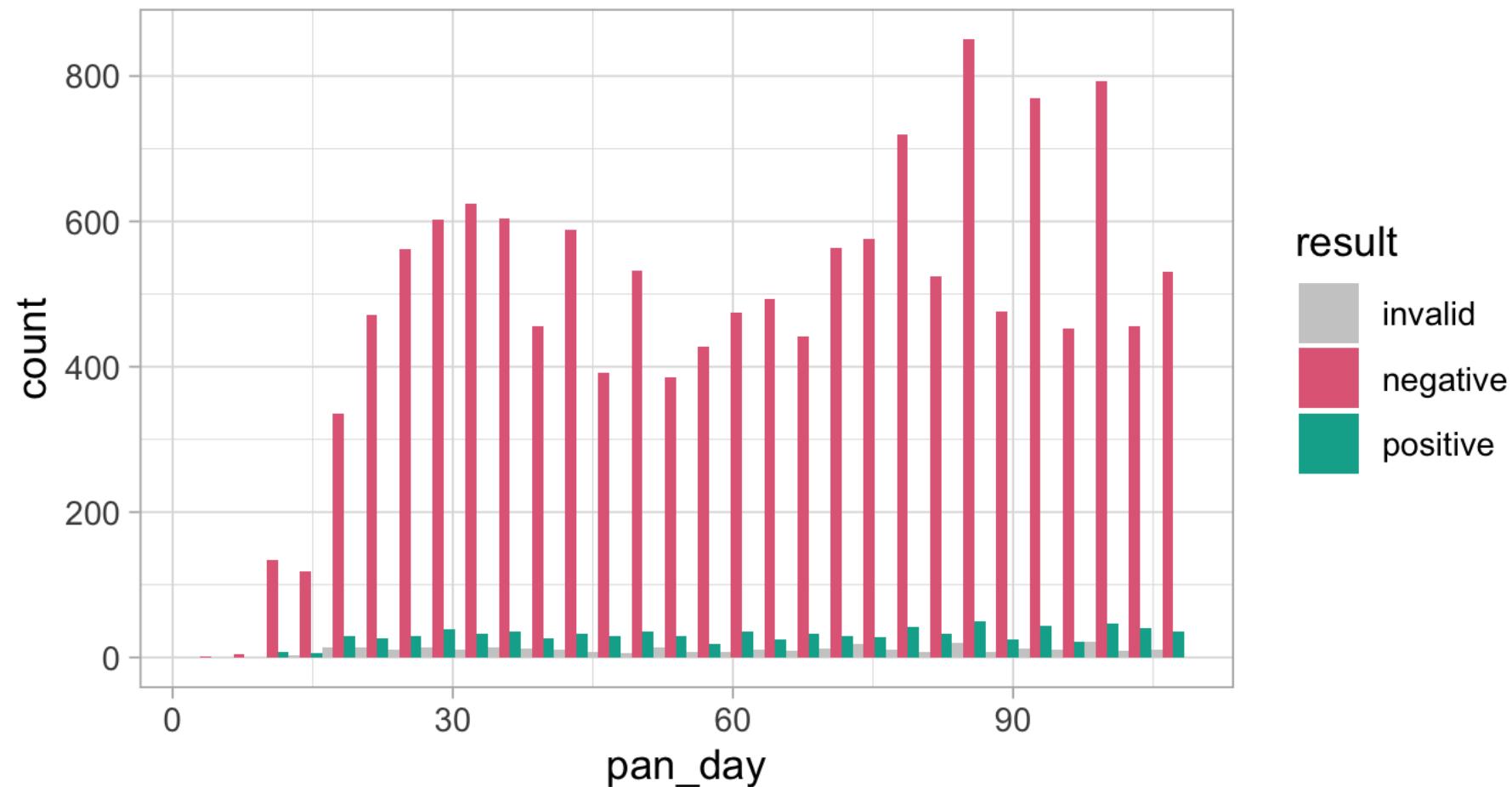
Themes

Visual appearance of non-data elements



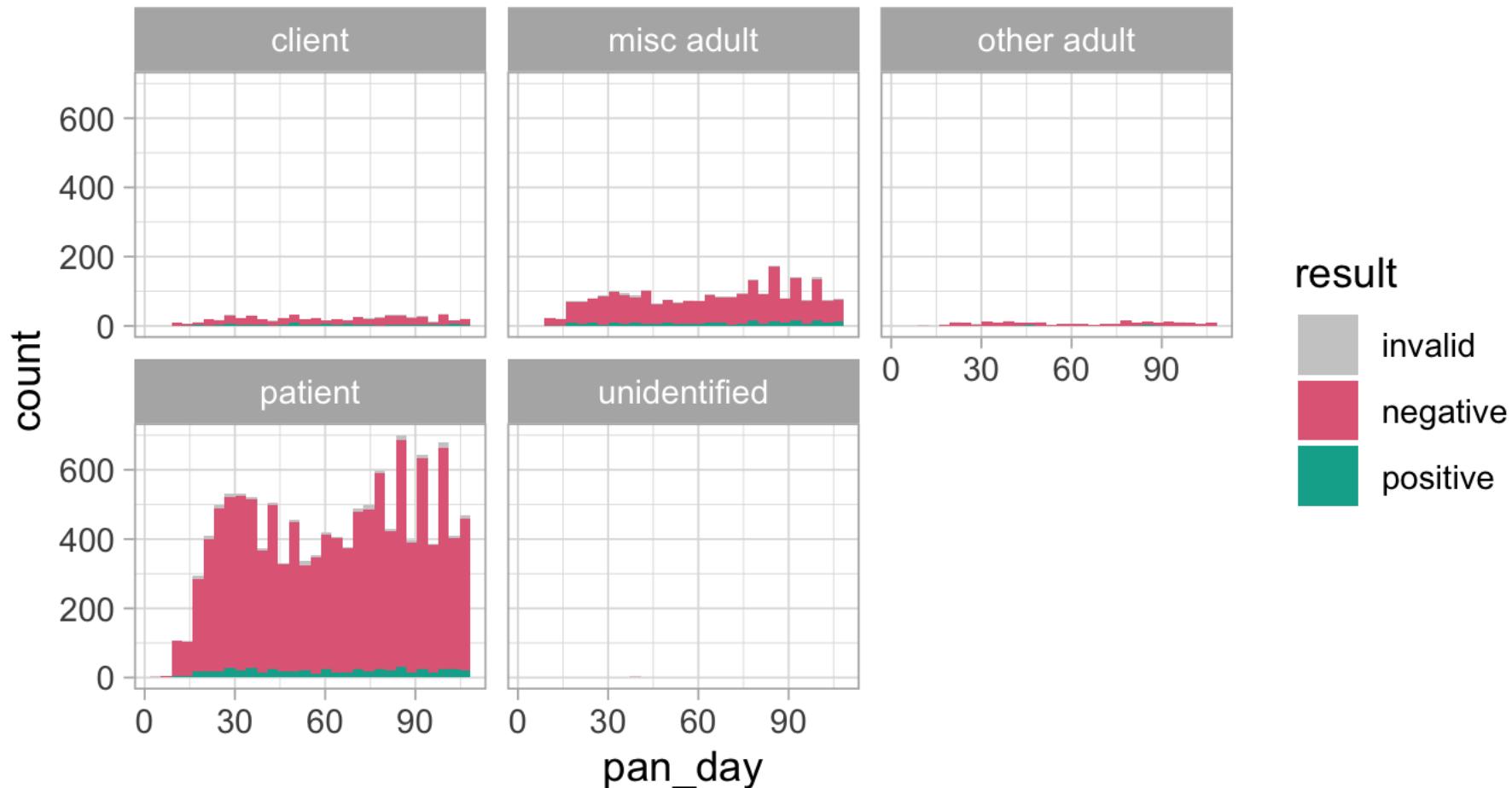
Scales

Customize color scales and other mappings

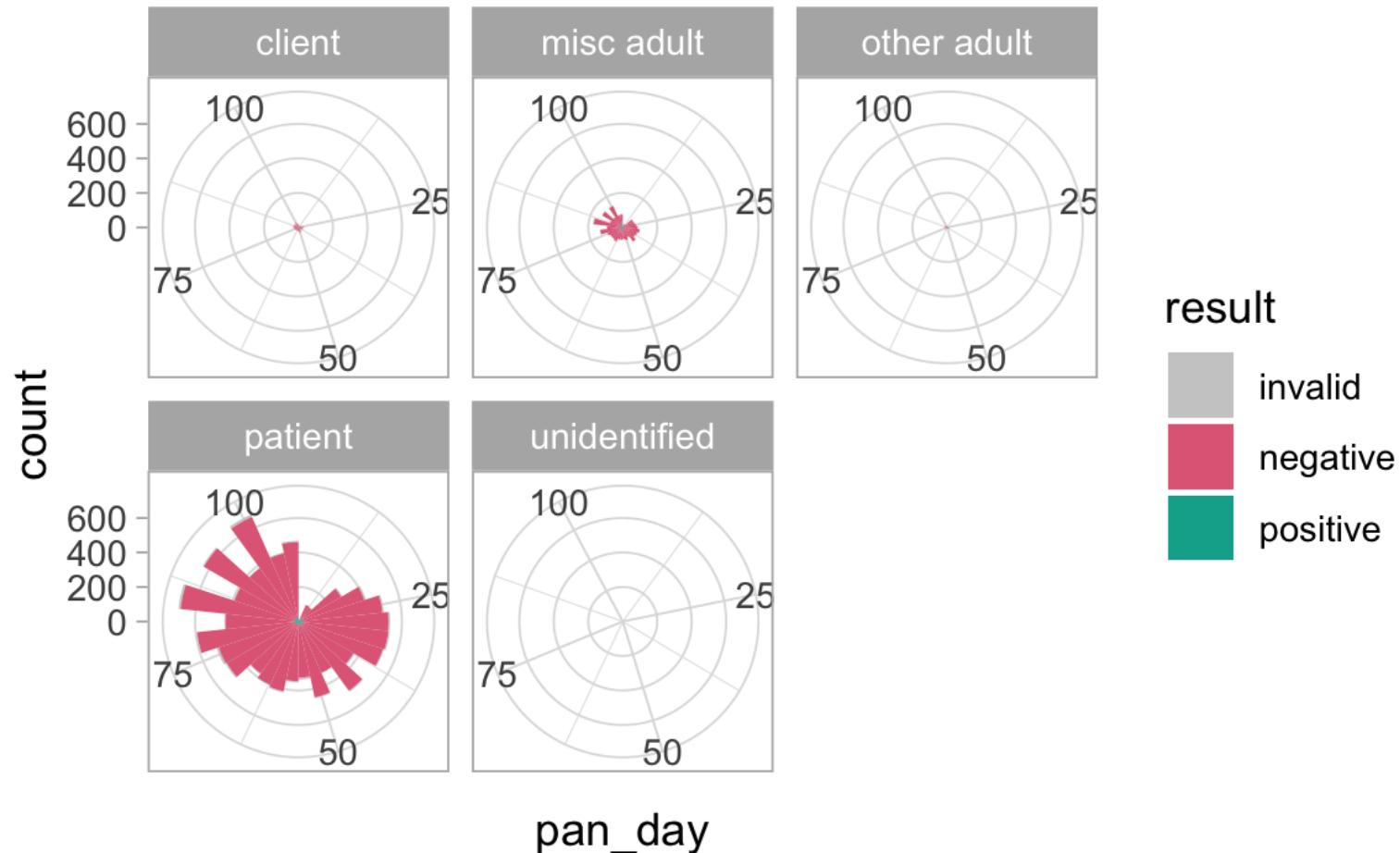


Facets

Subplots that display subsets of the data



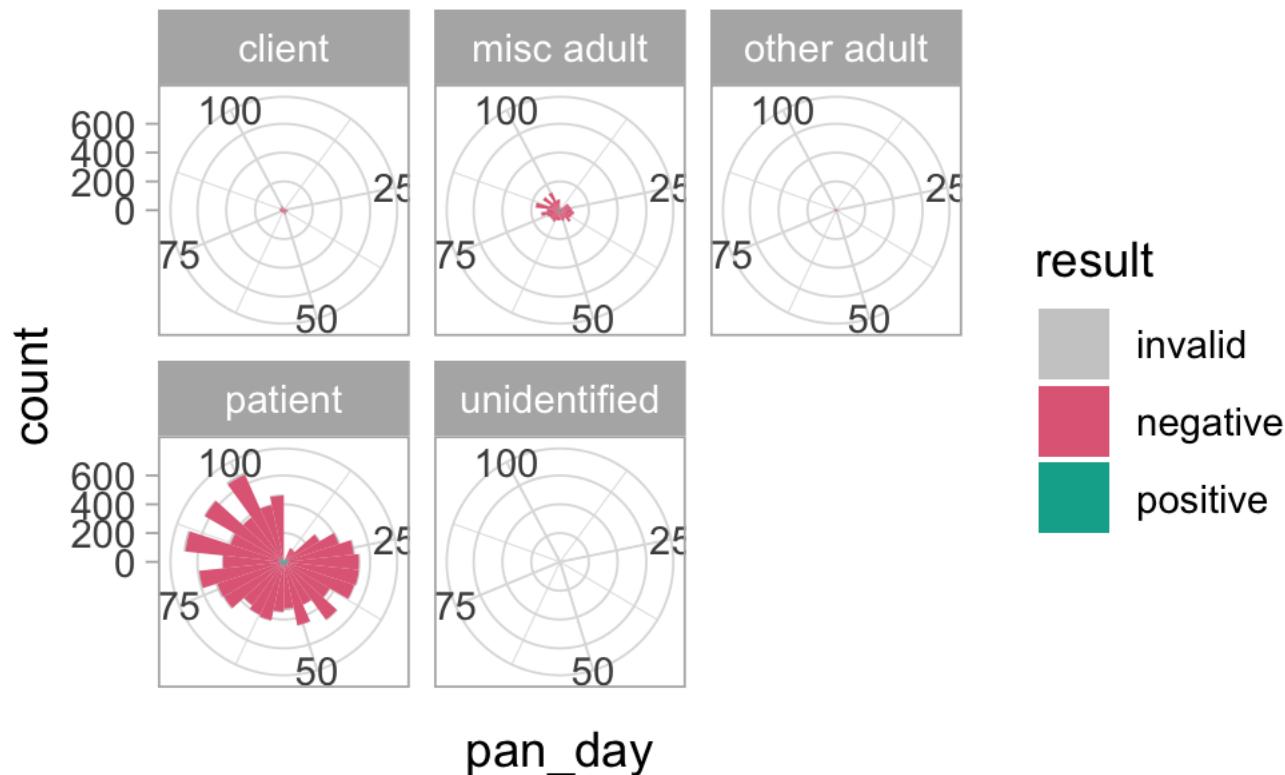
Coordinate systems



Titles and captions

COVID19 test volume

Displayed in polar coordinates, mostly to show off ggplot2



```
ggplot(data = data_frame) +  
  geom_function(mapping = aes(mappings)) +  
  theme_function +  
  scale_function +  
  facet_function +  
  coordinate_function +  
  ...
```

Required

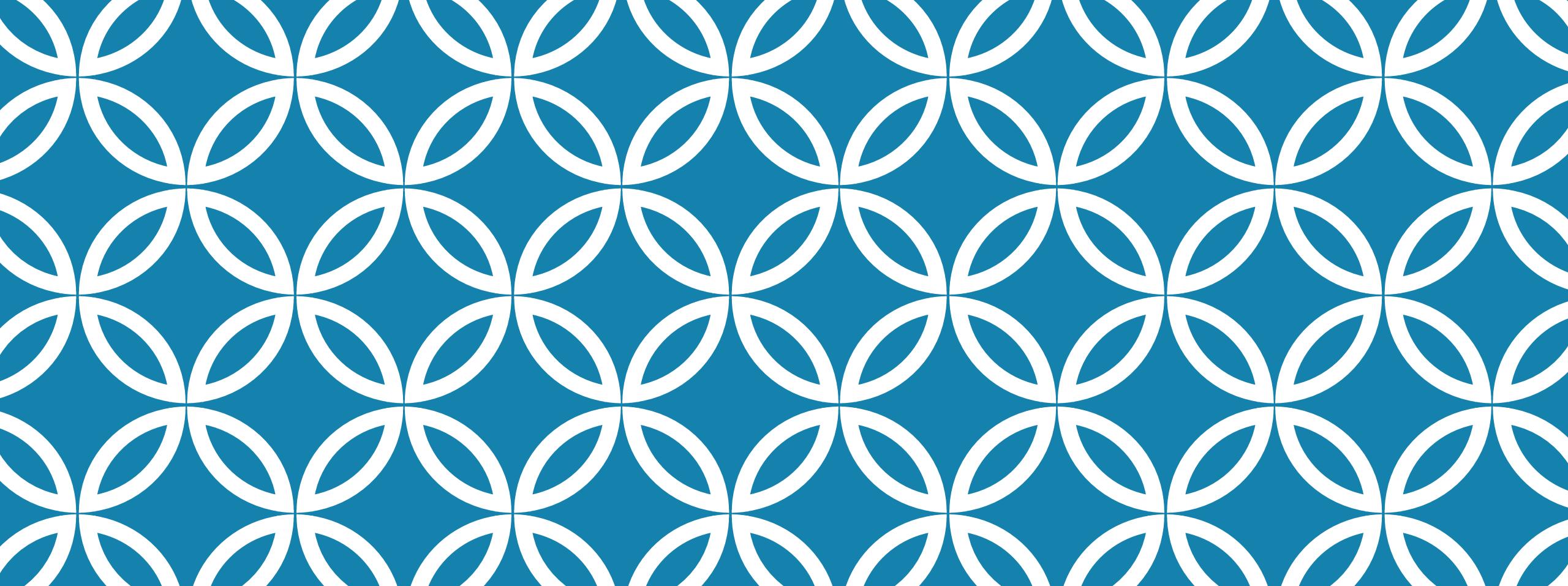
Optional

Goals

1. Appreciate the importance of visualization for understanding data
2. Learn how to use ggplot2 to visualize data

Objectives

1. Create a basic visualization using a simple template
2. Define “aesthetic mapping” and explain how aesthetic mappings relate variables of a data frame to features of graphic markings on a plot
3. Write the code to specify a type of plot and fine tune its appearance using “geom” functions
4. Explain how to add layers to a ggplot object to create complex and highly customized visualizations



Data Transformation

Session 4
Amrom Obstfeld
August 23, 2020

| August 16 2020 | Session | Instructor |
|---------------------|--|-------------------|
| 9:00 am - 9:30 am | Instructor Introductions, Introduction to technology | Amrom Obstfeld |
| 9:30 am - 10:15 am | Introduction to R and RStudio | Amrom Obstfeld |
| 10:30 pm - 11:15 am | Reproducible Reporting | Amrom Obstfeld |
| 11:30 am - 1:00 am | Data Visualization | Stephan Kadauke |
| August 23 2020 | | |
| 9:00 am - 10:30 pm | Data Transformation | Amrom Obstfeld |
| 10:45 am - 12:15 pm | Statistical Analysis | Dan Herman |
| 12:30 pm - 1:00 pm | Workshop Close out | Amrom Obstfeld |

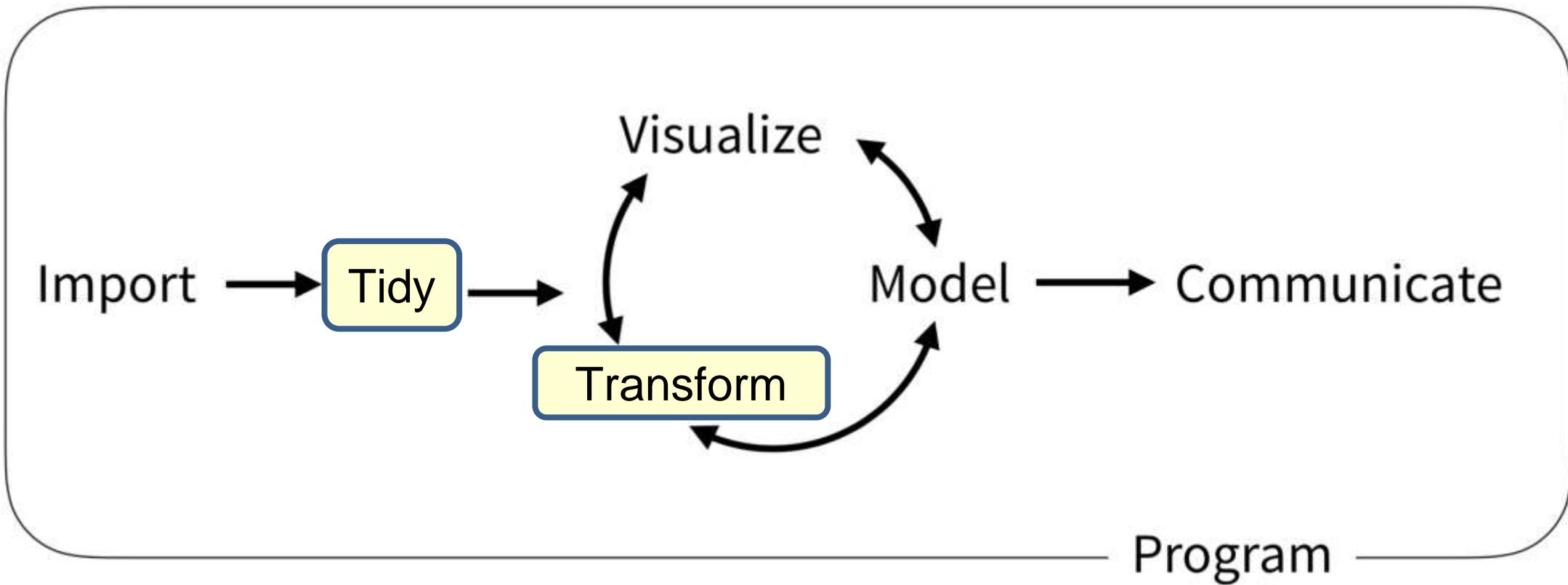
Goal

1. Learn how to use dplyr to transform data frames
2. Appreciate the role of piping in facilitating data transformation

Objectives

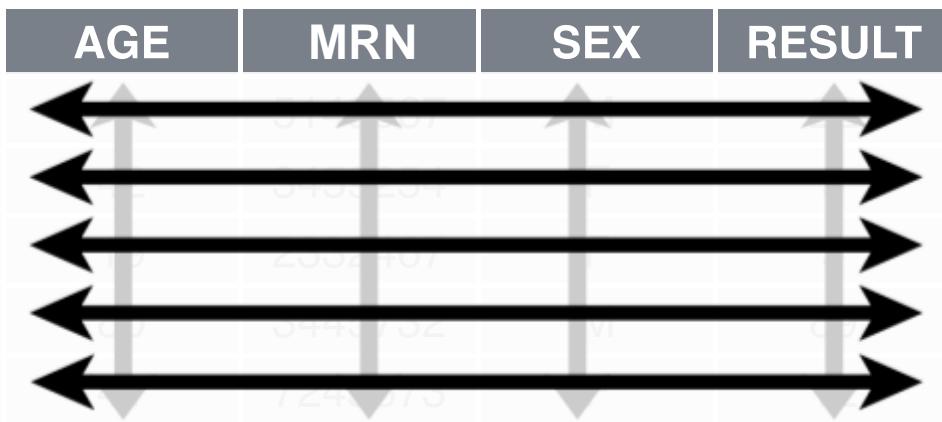
1. List the major forms of data transformation implemented in dplyr
2. Use code templates with dplyr functions to tidy a raw data set
3. Use the pipe operator to pass the output of one function as an input to the next function
4. Create new calculated columns not found in the original data frame

Typical Data Science Pipeline



What is a “Tidy” Data Frame

A data set is **tidy** if:



1. Each **variable** is in its own **column**
2. Each **observation** is in its own **row**
3. Each **value** is in its own **cell**

Your Turn 1

Open "04-Transform.Rmd"

Run the setup chunk

```
```{r setup}
library(tidyverse) # Provides functions used throughout this session

covid_testing <- read_csv("data/covid_testing.csv")
```
```



Pop Quiz

How can you confirm that you have successfully loaded the data file into Rstudio?

1. The code that imported the data did not yield an error
2. Code that references the `covid_testing` object runs without errors
3. The `covid_testing` object is present in the environment pane
4. All of the above

Transform Data with



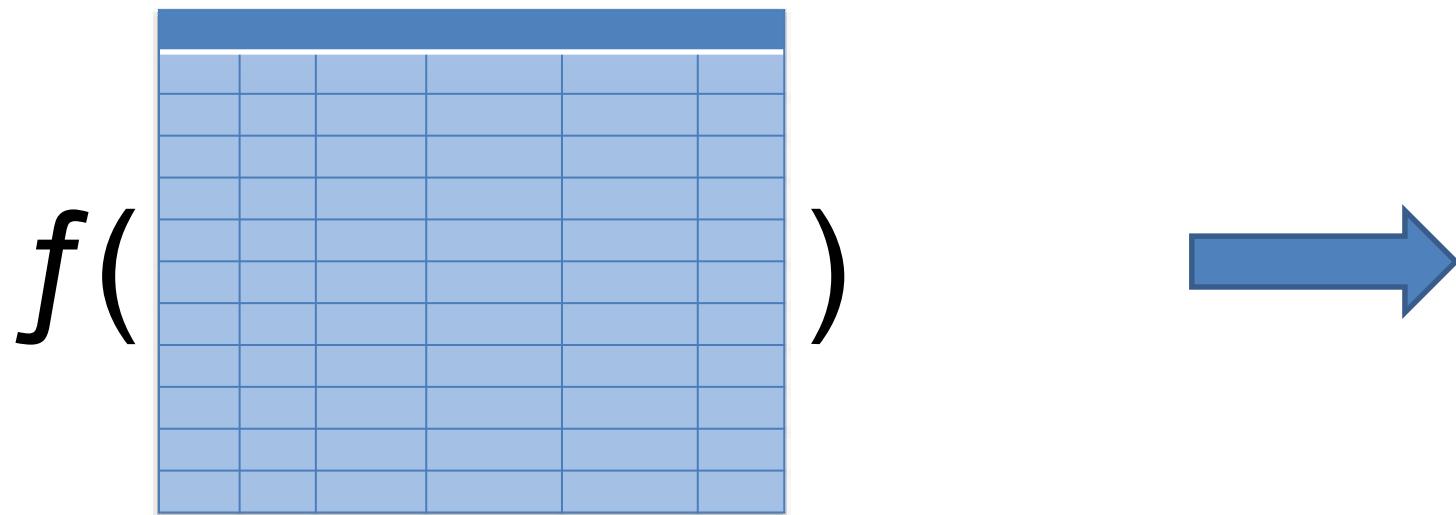
dplyr



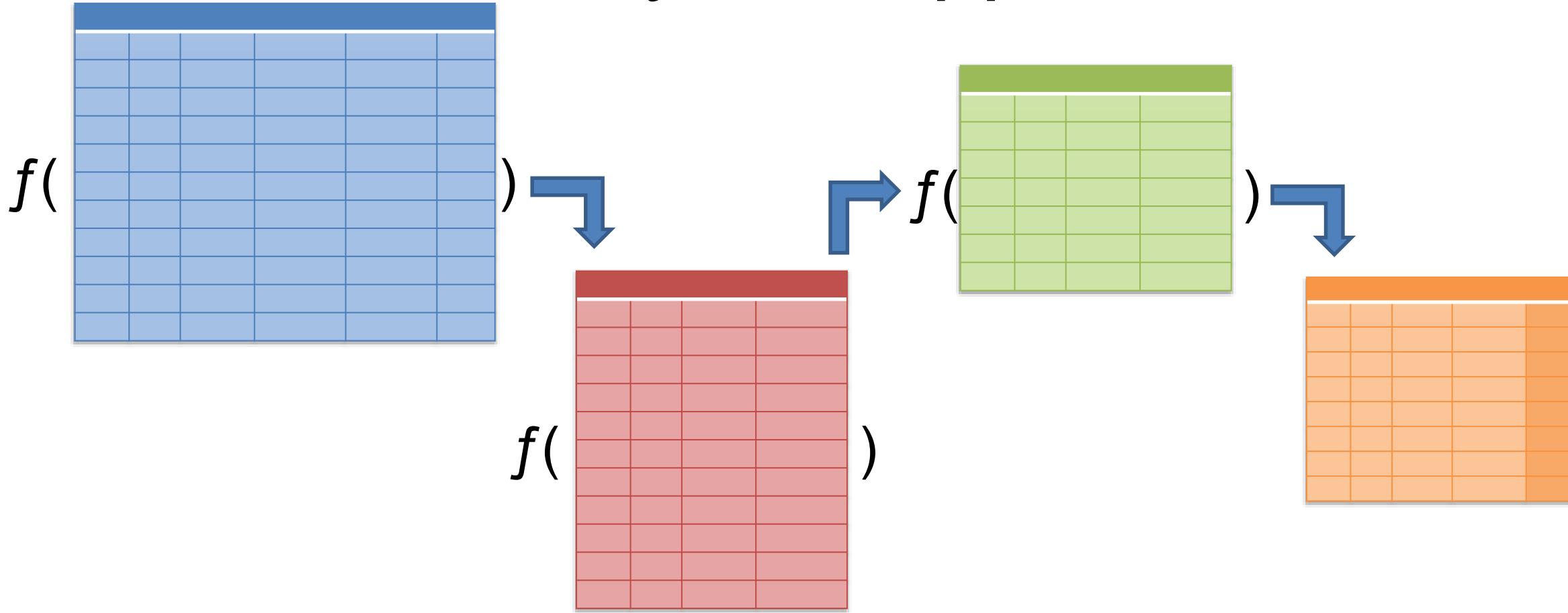
dplyr implements a *grammar* for transforming tabular data.



Analytical Approach



Analytical Approach



dplyr: a grammar for transforming data

- 1 Choose columns.** `select()`
- 2 Extract rows.** `filter()`
- 3 Derive new columns.** `mutate()`
- 4 Change the unit of analysis.** `summarize()`



Common syntax

Each function takes a data frame as its first argument and returns a data frame as its output.

```
function(data, ...)
```

dplyr
function

data frame to
transform

specific
arguments



Isolating data

select()

select()

Extract columns from a data frame

| | | | |
|------|------|------|------|
| Blue | Grey | Blue | Grey |
| Blue | Grey | Blue | Grey |
| Blue | Grey | Blue | Grey |
| Blue | Grey | Blue | Grey |
| Blue | Grey | Blue | Grey |



| | |
|------|--|
| Blue | |

= Number of rows
↓ Number of Columns



select()

Extract columns from a data frame

```
select(covid_testing, mrn, last_name)
```

dplyr
function

data frame to
transform

name(s) of columns
to extract
(or a select helper)



select()

Extract columns from a data frame **by name**

```
select(covid_testing, mrn, last_name)
```

covid_testing

| mrn | first_name | last_name | gender |
|---------|------------|------------|--------|
| 5000876 | sarella | stark | female |
| 5006017 | alester | stark | male |
| 5001412 | jhezane | westerling | female |
| 5000533 | penny | targaryen | female |
| ... | | | |



| mrn | last_name |
|---------|------------|
| 5000876 | stark |
| 5006017 | stark |
| 5001412 | westerling |
| 5000533 | targaryen |



select()

Extract columns from a data frame **by name**

```
select(covid_testing, -mrn, -last_name)
```

covid_testing

| mrn | first_name | last_name | gender |
|---------|------------|------------|--------|
| 5000876 | sarella | stark | female |
| 5006017 | alester | stark | male |
| 5001412 | jhezane | westerling | female |
| 5000533 | penny | targaryen | female |
| ... | | | |



| first_name | gender |
|------------|--------|
| sarella | female |
| alester | male |
| jhezane | female |
| penny | female |



select() helpers

Data Transformation with dplyr:: CHEAT SHEET

**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)**



Your Turn 2

- Alter the code to select just the `first_name` column from `covid_testing`
- If you have time, try to remove the `first_name` column

```
covid_testing_2 <- select(covid_testing, _____)
```

filter()

filter()

Extract rows that meet logical criteria

| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |



| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

↓ Number of rows
= Number of Columns



Common syntax

Each function takes a data frame as its first argument and returns a data frame as its output.

```
function(data, ...)
```

dplyr
function

data frame to
transform

specific
arguments



filter()

Extract rows that meet logical criteria

```
filter(data, ...)
```

data frame to
transform

one or more logical tests
(filter returns each row for
which the test is TRUE)

| | | | | |
|--|--|--|--|-------|
| | | | | FALSE |
| | | | | FALSE |
| | | | | TRUE |
| | | | | FALSE |
| | | | | TRUE |
| | | | | FALSE |



| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |



filter()

Extract rows that meet logical criteria

```
filter(data, column_name == criteria )
```

one or more logical tests
(filter returns each row for
which the test is TRUE)

| | | | | |
|--|--|--|--|-------|
| | | | | FALSE |
| | | | | FALSE |
| | | | | TRUE |
| | | | | FALSE |
| | | | | TRUE |
| | | | | FALSE |



| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |



filter()

Extract rows that meet logical criteria

```
filter(covid_testing, mrn==5000083)
```

| | mrn | first_name | last_name |
|-------|---------|------------|------------|
| FALSE | 5000876 | sarella | stark |
| FALSE | 5006017 | alester | stark |
| FALSE | 5001412 | jhezane | westerling |
| TRUE | 5000083 | lollys | clegane |



| | mrn | first_name | last_name |
|--|---------|------------|-----------|
| | 5000083 | lollys | clegane |



filter()

Extract rows that meet logical criteria

```
filter(covid_testing, mrn==5000083)
```

| mrn | first_name | last_name |
|---------|------------|------------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |
| 5001412 | jhezane | westerling |
| 5000083 | lollys | clegane |

= sets

(returns nothing)

== tests if equal

(returns TRUE or FALSE)



filter()

Values coded as character strings must be surrounded by quotes

Extract rows that meet logical criteria.

```
filter(covid_testing, last_name=="stark")
```

| mrn | first_name | last_name | |
|---------|------------|------------|-------|
| 5000876 | sarella | stark | TRUE |
| 5006017 | alester | stark | TRUE |
| 5001412 | jhezane | westerling | FALSE |
| 5000083 | lollys | clegane | FALSE |



| mrn | first_name | last_name |
|---------|------------|-----------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |



filter()

Extract rows that meet logical criteria

```
filter(data, ... )
```

data frame to
transform

one or more logical tests
(filter returns each row for
which the test is TRUE)



Logical tests

| | |
|------------------------|--------------------------|
| <code>x < y</code> | Less than |
| <code>x > y</code> | Greater than |
| <code>x == y</code> | Equal to |
| <code>x <= y</code> | Less than or equal to |
| <code>x >= y</code> | Greater than or equal to |
| <code>x != y</code> | Not equal to |
| <code>x %in% y</code> | Group membership |
| <code>is.na(x)</code> | Is NA |
| <code>!is.na(x)</code> | Is not NA |



Pop Quiz

What is the result?

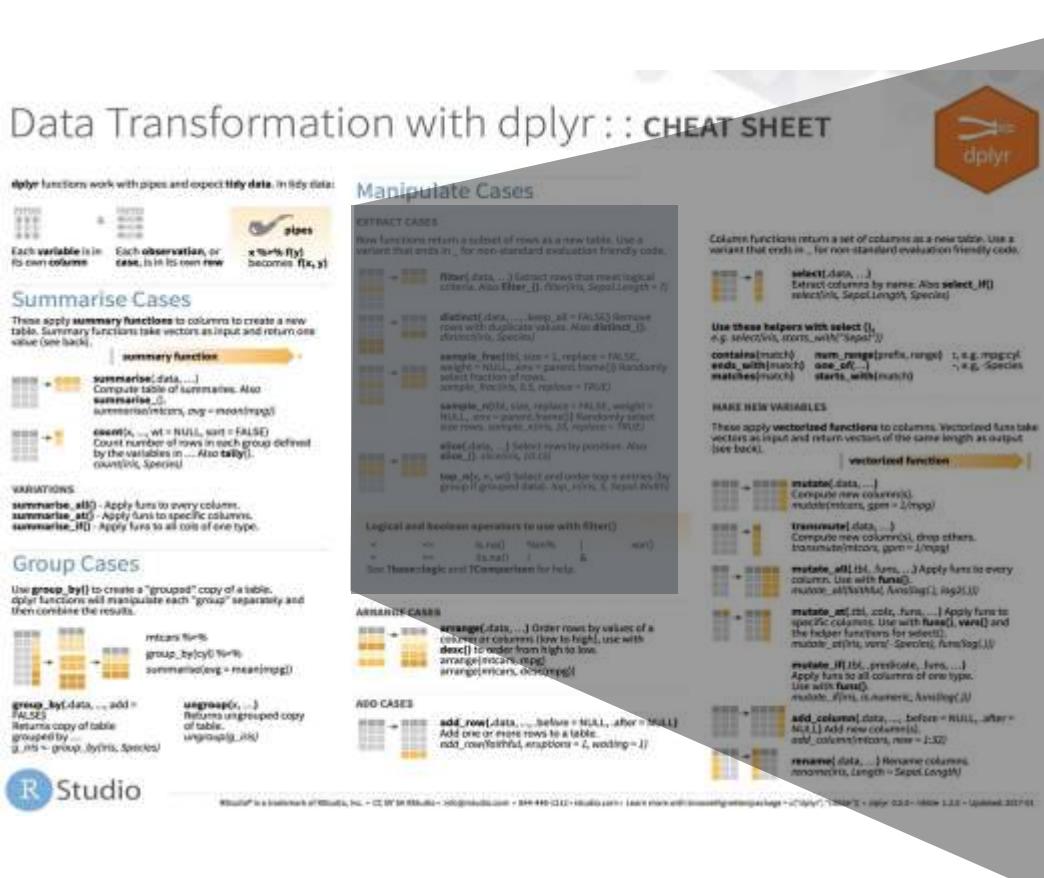
1 == 1

Pop Quiz

What is the result?

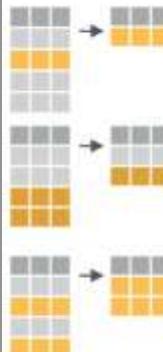
$$3 \neq 1$$

filter() variants

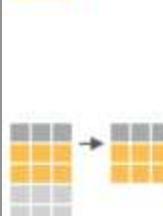


EXTRACT CASES

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



filter(.data, ...) Extract rows that meet logical criteria. `filter(iris, Sepal.Length > 7)`



distinct(.data, ..., .keep_all = FALSE) Remove rows with duplicate values.
distinct(iris, Species)

```
sample_frac(tbl, size = 1, replace = FALSE,  
weight = NULL, .env = parent.frame()) Randomly  
select fraction of rows.  
sample_frac(iris, 0.5, replace = TRUE)
```

sample_n(tbl, size, replace = FALSE, weight = NULL, .env = parent.frame()) Randomly select size rows. *sample_n(iris, 10, replace = TRUE)*

slice(.data, ...) Select rows by position.
`slice(iris, 10:15)`

top_n(x, n, wt) Select and order top n entries (by group if grouped data). *top_n(iris, 5, Sepal.Width)*

Logical and boolean operators to use with filter()

| | | | | | |
|---|--------|-----------------------|-------------------|--------------------|--------------------|
| < | \leq | <code>is.na()</code> | <code>%in%</code> | | <code>xor()</code> |
| > | \geq | <code>!is.na()</code> | <code>!</code> | <code>&</code> | |

See [?base::logic](#) and [?Comparison](#) for help.



Your Turn 3

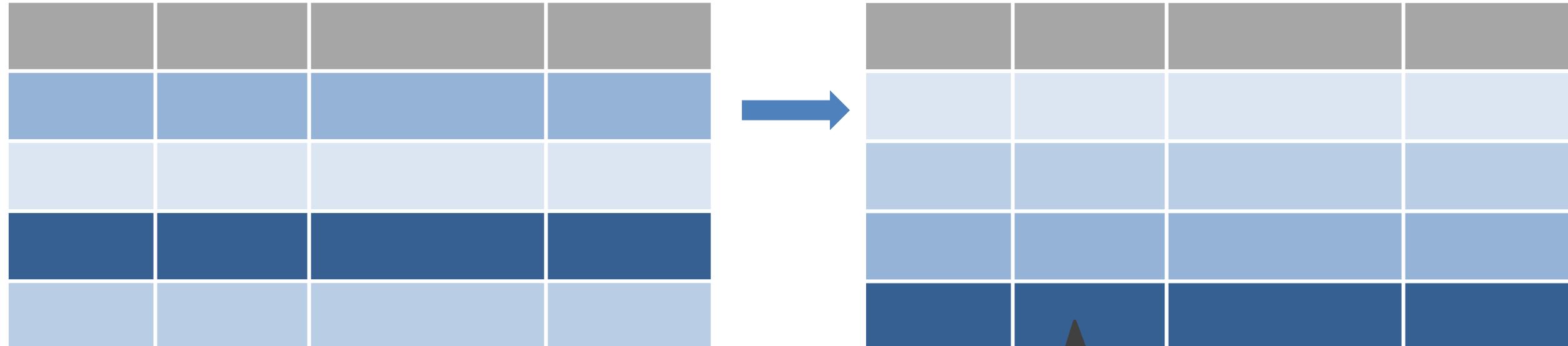
Use filter() with the logical operators to find:

- Every test for patients **over age 80**
- All of the covid testing where the demographic group (demo_group) is **equal to "client"**
- CHALLENGE:
 - All of the covid testing where the patient class (patient_class) **is NA** [Hint: See slide titled "Logical Tests"]

arrange()

arrange()

Order rows by values in a column



- = Number of rows
- = Number of Columns



arrange()

Order rows by values in a column

```
arrange(data, ... )
```

data frame to
transform

name(s) of columns to
arrange by



arrange()

Order rows by values in a column

```
arrange(covid_testing, first_name)
```

| mrn | first_name | last_name |
|---------|------------|-----------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |
| 5001412 | jhezane | targaryen |
| 5000533 | penny | targaryen |



| mrn | first_name | last_name |
|---------|------------|-----------|
| 5006017 | alester | stark |
| 5001412 | jhezane | targaryen |
| 5000533 | penny | targaryen |
| 5000876 | sarella | stark |



arrange()

Order rows by values in a column

```
arrange(covid_testing, desc(mrn))
```

| mrn | first_name | last_name |
|---------|------------|-----------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |
| 5001412 | jhezane | targaryen |
| 5000533 | penny | targaryen |



| mrn | first_name | last_name |
|---------|------------|-----------|
| 5006017 | alester | stark |
| 5001412 | jhezane | targaryen |
| 5000876 | sarella | stark |
| 5000533 | penny | targaryen |



Your Turn 4

The column `ct_value` contains the cycle threshold (Ct) for the real-time PCR that generated the final result.

How might you use `arrange()` to determine the highest and lowest Ct result in the dataset?

Pop Quiz

The default behavior of `arrange()` is to order from lower to higher values.

When might `arrange()` place "1000" before "50"?

%>%

Data Analysis Steps

```
day_10 <- filter(covid_testing, pan_day <= 10)  
day_10 <- select(day_10, clinic_name)  
day_10 <- arrange(day_10 , clinic_name)
```

1. Filter tests to those on pandemic day less than 10
2. Select the column that contains ordering location
3. Arrange those columns by location



Data Analysis Steps

```
day_10 <- arrange(  
  select(  
    filter(  
      covid_testing,  
      pan_day <= 10  
    ),  
    clinic_name  
  ),  
  clinic_name  
)
```



The Pipe Operator %>%

Passes result on left into first argument of function on right.



```
covid_testing %>% filter(_____, pan_day <= 10)
```

```
filter(covid_tesing, pan_day <= 10)  
covid_tesing %>% filter(pan_day <= 10)
```



Data Analysis Steps

```
day_10 <- arrange(  
  select(  
    filter(  
      covid_testing,  
      pan_day <= 10  
    ),  
    clinic_name  
  ),  
  clinic_name  
)
```



Data Analysis Steps

```
covid_testing %>%  
  filter(pan_day <= 10) %>%  
  select(clinic_name) %>%  
  arrange(clinic_name)
```



Shortcut to type %>%

Cmd + **Shift** + **M** (Mac)

Ctrl + **Shift** + **M** (Windows)



Scene

The PICU would like a word with you because of a recent incident involving a delay in results for a patient who required a AGP

They had to wait over 10 hours before the procedure could begin

You decide to investigate... WITH DATA



Question

Was the TAT so high in this situation because:

- A. Your lab has long TAT's and you need to find a way to get the results to certain priority areas faster
- B. Something was missed in this one case and the workflow just needs to be reviewed and tweaked so things aren't missed



Your Turn 5

Use `%>%` to write a sequence of three functions that:

1. Filters to tests from the clinic (`clinic_name`) of "picu"
2. Selects the column with the receive to verify turnaround time (`rec_ver_tat`) as well as the day from start of the pandemic (`pan_day`)
3. Arrange the `pan_day` from highest to lowest

Using `<-`, assign the result to a new variable, call it whatever you want.

Isolating data



Extract variables with `select()`



Extract rows with `filter()`



Arrange rows, with `arrange()`.



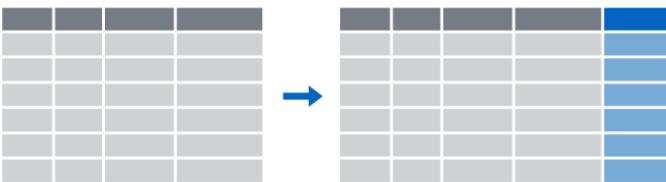
Deriving Data

**What is the average
and SD in total
turnaround time by
clinic?**

Breaking down the analytical question

1. Total TAT for each test
2. Group tests by clinic
3. Calculate average and SD for each clinic

Deriving data



Make new variables with **mutate()**



Make summaries of data with
summarize()



mutate()

mutate()

Creating new calculated columns



= Number of rows
↑ Number of Columns



mutate()

Creating new calculated columns

```
covid_testing %>%  
  mutate(new_column = calculation)
```

name for new column

equals

function whose results will populate columns



mutate()

Creating new calculated columns

```
covid_testing %>%  
  mutate(c_r_tat_mins = col_rec_tat * 60)
```

| mrn | col_rec_tat | rec_ver_tat |
|---------|-------------|-------------|
| 5000876 | 29.5 | 11.5 |
| 5006017 | 3.6 | 5 |
| 5001412 | 1.4 | 5.2 |
| 5000533 | 2.3 | 5.8 |



| mrn | col_rec_tat | rec_ver_tat | c_r_tat_mins |
|---------|-------------|-------------|--------------|
| 5000876 | 29.5 | 11.5 | 1770 |
| 5006017 | 3.6 | 5 | 216 |
| 5001412 | 1.4 | 5.2 | 84 |
| 5000533 | 2.3 | 5.8 | 138 |

Your Turn 6

Create a new column using the `mutate()` function that contains the total TAT (sum of `col_rec_tat` and `rec_ver_tat`)

mutate()

Replacing columns

```
covid_testing %>%  
  mutate(mrn = as.character(mrn))
```

Function to "coerce" one type of data into another type of data

| mrn
<dbl> | first_name
<chr> | last_name
<chr> |
|--------------|---------------------|--------------------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |
| 5001412 | jhezane | westerling |
| 5000533 | penny | targaryen |



| mrn
<chr> | first_name
<chr> | last_name
<chr> |
|--------------|---------------------|--------------------|
| 5000876 | sarella | stark |
| 5006017 | alester | stark |
| 5001412 | jhezane | westerling |
| 5000533 | penny | targaryen |



Functions to use in mutate()

| |
|---|
| <h1>Vector Functions</h1> |
| <h2>TO USE WITH MUTATE ()</h2> |
| <p>mutate() and transmute() apply vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.</p> |
| <h3>vectorized function</h3> |
| <h4>OFFSETS</h4> |
| <p>dplyr::lag() - Offset elements by 1
dplyr::lead() - Offset elements by -1</p> |
| <h4>CUMULATIVE AGGREGATES</h4> |
| <p>dplyr::cumall() - Cumulative all()
dplyr::cumany() - Cumulative any()
 cummax() - Cumulative max()
dplyr::cummean() - Cumulative mean()
 cummin() - Cumulative min()
 cumprod() - Cumulative prod()
 cumsum() - Cumulative sum()</p> |
| <h4>RANKINGS</h4> |
| <p>dplyr::cume_dist() - Proportion of all values <=
dplyr::dense_rank() - rank with ties = min, no gaps
dplyr::min_rank() - rank with ties = min
dplyr::ntile() - bins into n bins
dplyr::percent_rank() - min_rank scaled to [0,1]
dplyr::row_number() - rank with ties = "first"</p> |
| <h4>MATH</h4> |
| <p>+, -, *, /, ^, %%, %% - arithmetic ops
log(), log2(), log10() - logs
<, <=, >, >=, !=, == - logical comparisons</p> |
| <p>dplyr::between() - x >= left & x <= right
dplyr::near() - safe == for floating point numbers</p> |
| <h4>MISC</h4> |
| <p>dplyr::case_when() - multi-case if_else()
dplyr::coalesce() - first non-NA values by element across a set of vectors
dplyr::if_else() - element-wise if() + else()
dplyr::na_if() - replace specific values with NA
 pmax() - element-wise max()
 pmin() - element-wise min()
dplyr::recode() - Vectorized switch()
dplyr::recode_factor() - Vectorized switch() for factors</p> |

Vector Functions

TO USE WITH MUTATE()

- `mutate()` and `transmute()` apply vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.

vectorized function

Summary Functions

TO USE WITH SUMMARISE()

- `summarise()` applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

summary function

Combine Tables

COMBINE VARIABLES

`bind_rows()` + `bind_cols()` =

Use `bind_collapse()` to paste tables beside each other as they are.

COMBINE CASES

`bind_rows()` + `bind_col()` =

Use `bind_new()` to paste tables below each other as they are.

OTHER FUNCTIONS

COUNTS

- `n()`: number of values in row
- `n_distinct()`: # of unique values
- `name_if(is.na)`: if not NA

LOCATION

- `mean()`: mean, also `mean(is.na())`
- `median()`: median

LOGICALS

- `isTRUE()`: Properties of TRUE
- `isTRUE(#)`: # of TRUE's

POSITION/NUMBER

- `plyr::first()`: first value
- `plyr::last()`: last value
- `plyr::with()`: value in nth location of vector

NAME

- `quantile()`: nth quantile
- `min()`: minimum value
- `max()`: maximum value

SPREAD

- `sq()`: Inter-Quantile Range
- `mad()`: median absolute deviation
- `sd()`: standard deviation
- `var()`: variance

Row Names

This data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.

`rownames_to_column()`
 Move row names into col.
`g ~ rownames(.), to = column_name, vrt = "T"`

`column_to_rownames()`
 Move col in row names.
`inplace = TRUE, row.names = "C"`

`rowwise_(, remove_rownames[])`
 Also `rowwise_(, remove_rownames[])`

Combine Tables

COMBINE VARIABLES

`bind_rows()` + `bind_cols()` =

Use `bind_collapse()` to paste tables beside each other as they are.

COMBINE CASES

`bind_rows()` + `bind_col()` =

Use `bind_new()` to paste tables below each other as they are.

OTHER FUNCTIONS

INTERSECT

`intersect(x, y, ...)`
 Rows that appear in both x and y.

SIMPLY

`setdiff(x, y, ...)`
 Rows that appear in x but not y.

UNION

`union(x, y, ...)`
 Rows that appear in x or y.
`unique(..., drop = TRUE)`
 (Duplicates removed). `unique(..., drop = TRUE)` retains duplicates.

EQUALITY

`equate(x, y)`
 Use `equate()` to test whether two data sets contain the exact same rows (in any order).

EXTRACT ROWS

`filter(x, ..., .by = y)`
 Use `filter()` to filter one table against the rows of another.

SELECT COLUMNS

`select(x, ..., .by = y, ...)`
 Return rows of x that have a match in y.
`useful = TRUE` TO SEE WHAT WILL BE JOINED

JOIN

`anti_join(x, y, ...)`
 Return rows of x that do not have a match in y. USEFUL TO SEE WHAT WILL NOT BE JOINED.

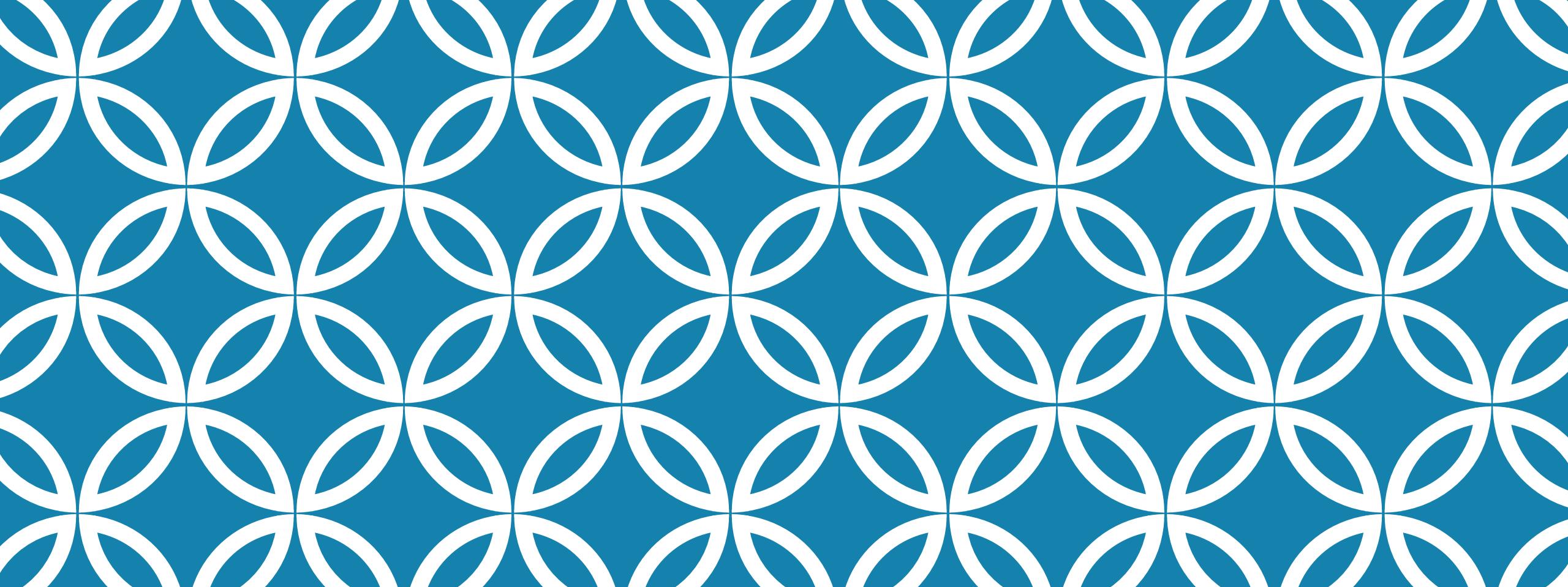


Goal

1. Learn how to use dplyr to transform data frames
2. Appreciate the role of piping in facilitating data transformation

Objectives

1. List the major forms of data transformation implemented in dplyr
2. Use code templates containing dplyr functions to tidy a raw data set
3. Use the pipe operator to pass the output of one function as an input to the next function
4. Create new calculated columns not found in the original data frame



Summary and Stats

Session 5
Dan Herman

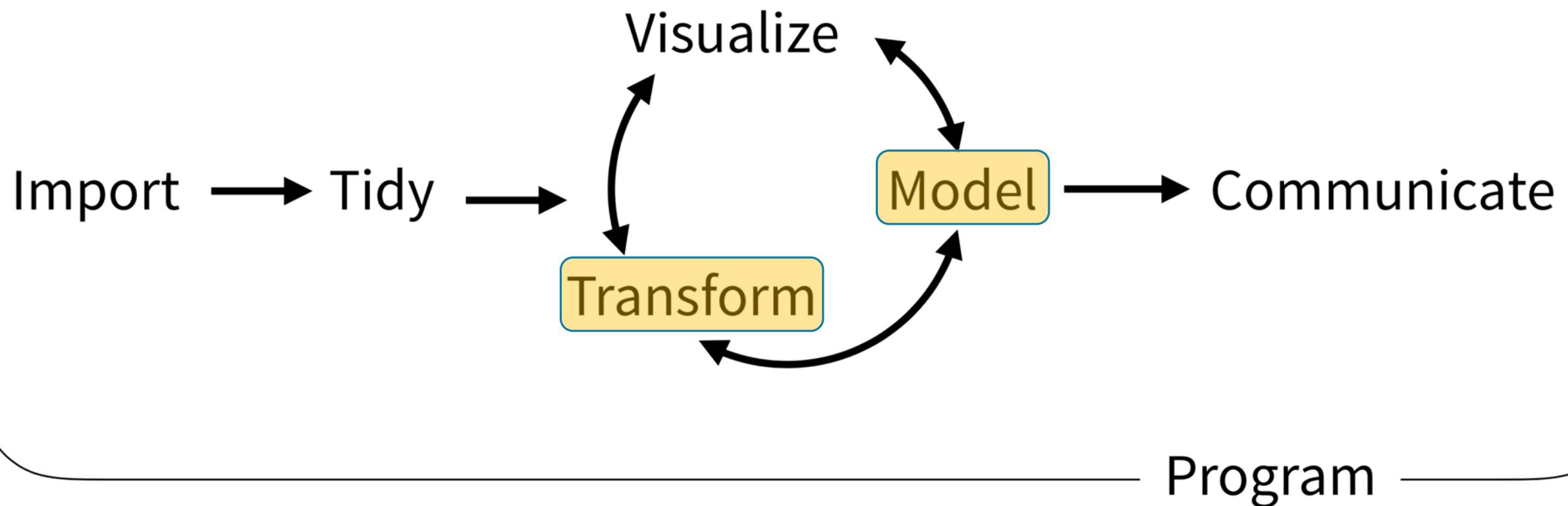
Goals

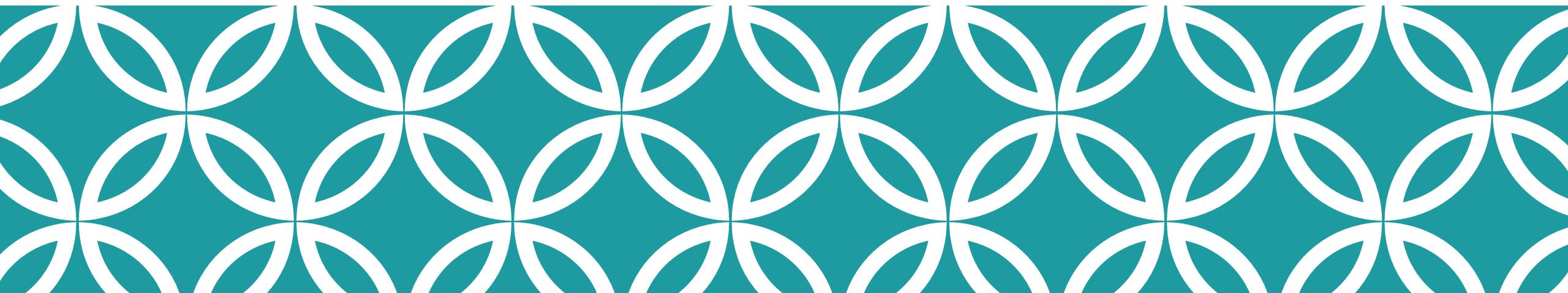
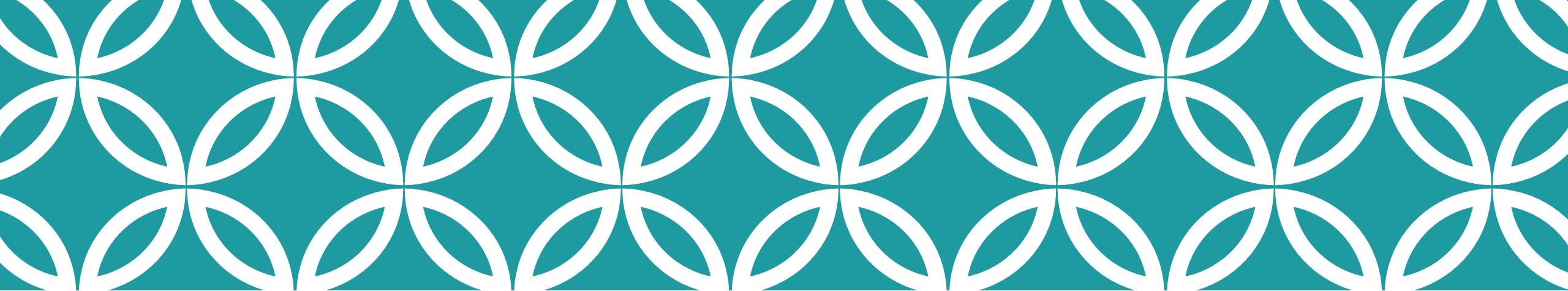
1. Learn how to summarize data
2. Appreciate how to perform hypothesis testing

Objectives

1. Calculate a summary statistic for a variable using `summarize`
2. Calculate summary statistics for a variable separately for a group of observations, using `group_by` and `summarize`
3. Describe the process of performing a simple test for association

Typical Data Science Pipeline

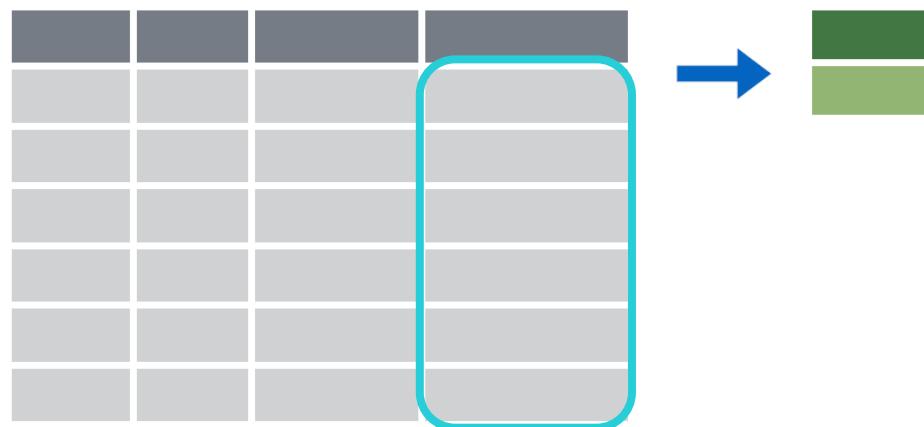




Summarize()

summarize()

- Make summaries of your data



summarize()

- Make summaries of your data

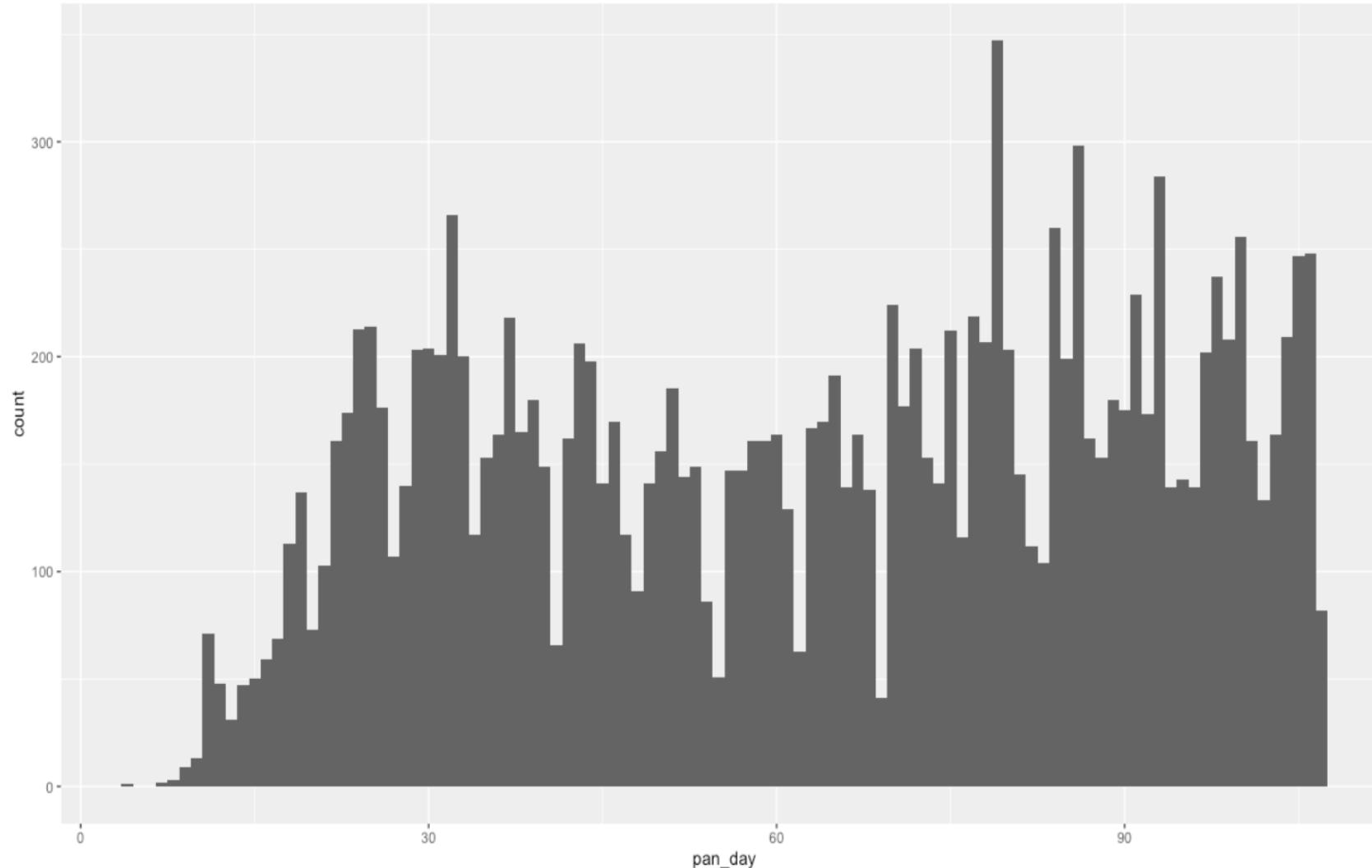
```
covid_testing %>%  
  summarize(new_variable = calculation)
```

name for new
variable

Value or
function



Q: How many tests are ordered per day?



summarize()

- Make summaries of your data

```
covid_testing %>%  
  select(mrn, pan_day) %>%  
  head(4) %>%  
  
  summarize(order_count = n())
```

function that returns
number of observations

| mrn | pan_day |
|---------|---------|
| 5001412 | 4 |
| 5000533 | 7 |
| 5009134 | 7 |
| 5008518 | 8 |



| order_count |
|-------------|
| 4 |

summarize()

- Make summaries of your data

```
covid_testing %>%  
  select(mrn, pan_day) %>%  
  head(4) %>%  
  
  summarize(order_count = n(),  
            day_count = n_distinct(pan_day))
```

function that returns
number of distinct values

| mrn | pan_day |
|---------|---------|
| 5001412 | 4 |
| 5000533 | 7 |
| 5009134 | 7 |
| 5008518 | 8 |



| order_count | day_count |
|-------------|-----------|
| 4 | 3 |

Your Turn 1

- Open “05 - Stats.Rmd”
- Run the setup chunk
- Fill-in gaps to calculate:
 - a) Mean count of orders per `pan_day`
 - b) Mean count of orders per clinic



Vector Functions

TO USE WITH MUTATE()

COUNTS

`dplyr::n()` - number of values/rows
`dplyr::n_distinct()` - # of uniques
`sum(!is.na())` - # of non-NA's

LOCATION

`mean()` - mean, also `mean(!is.na())`
`median()` - median

LOGICALS

`mean()` - Proportion of TRUE's
`sum()` - # of TRUE's

POSITION/ORDER

`dplyr::first()` - first value
`dplyr::last()` - last value
`dplyr::nth()` - value in nth location of vector

RANK

`quantile()` - nth quantile
`min()` - minimum value
`max()` - maximum value

SPREAD

`IQR()` - Inter-Quartile Range
`mad()` - median absolute deviation
`sd()` - standard deviation
`var()` - variance

Summary Functions

TO USE WITH SUMMARISE()

`summarise()` applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

summary function

COUNTS

`dplyr::n()` - number of values/rows
`dplyr::n_distinct()` - # of uniques
`sum(!is.na())` - # of non-NA's

LOCATION

`mean()` - mean, also `mean(!is.na())`
`median()` - median

LOGICALS

`mean()` - Proportion of TRUE's
`sum()` - # of TRUE's

POSITION/ORDER

`dplyr::first()` - first value
`dplyr::last()` - last value
`dplyr::nth()` - value in nth location of vector

RANK

`quantile()` - nth quantile
`min()` - minimum value
`max()` - maximum value

SPREAD

`IQR()` - Inter-Quartile Range
`mad()` - median absolute deviation
`sd()` - standard deviation
`var()` - variance

Row Names

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.

`rownames_to_column()`
 Move row names into col.
`a <- rownames_to_column(iris, var = "C")`

`column_to_rownames()`
 Move col in row names.
`column_to_rownames(a, var = "C")`

Also `has_rownames()`, `remove_rownames()`

Combine Tables

COMBINE VARIABLES

| x | y | = |
|-------|-------|-------------|
| A B C | A B D | A B C A B D |
| a t 1 | a t 3 | a t 1 a t 3 |
| b u 2 | b u 2 | b u 2 b u 2 |
| c v 3 | d w 1 | c v 3 d w 1 |

Use `bind_cols()` to paste tables beside each other as they are.

`bind_cols(...)` Returns tables placed side by side as a single table.
 BE SURE THAT ROWS ALIGN.

Use a "Mutating Join" to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

`left_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"), ...)`
 Join matching values from y to x.

`right_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"), ...)`
 Join matching values from x to y.

`inner_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"), ...)`
 Join data. Retain only rows with matches.

`full_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"), ...)`
 Join data. Retain all values, all rows.

Use `by = c("col1", "col2")` to specify the column(s) to match on.
`left_join(x, y, by = "A")`

Use a named vector, `by = c("col1" = "col2")`, to match on columns with different names in each data set.
`left_join(x, y, by = c("C" = "D"))`

Use `suffix` to specify suffix to give to duplicate column names.
`left_join(x, y, by = c("C" = "D"), suffix = c("1", "2"))`

COMBINE CASES

| x | y | = |
|-------|-------|-------------|
| A B C | A C V | A B C A C V |
| a t 1 | d w 4 | a t 1 d w 4 |

Use `bind_rows()` to paste tables below each other as they are.

`bind_rows(..., .id = NULL)`
 Returns tables one on top of the other as a single table. Set `.id` to a column name to add a column of the original table names (as pictured)

`intersect(x, y, ...)`
 Rows that appear in both x and y.

`setdiff(x, y, ...)`
 Rows that appear in x but not y.

`union(x, y, ...)`
 Rows that appear in x or y.
 (Duplicates removed). `union_all()` retains duplicates.

Use `setequal()` to test whether two data sets contain the exact same rows (in any order).

EXTRACT ROWS

| x | y | = |
|-------|-------|-------------|
| A B C | A B D | A B C A B D |
| a t 1 | a t 3 | a t 1 a t 3 |

Use a "Filtering Join" to filter one table against the rows of another.

`semi_join(x, y, by = NULL, ...)`
 Return rows of x that have a match in y.
 USEFUL TO SEE WHAT WILL BE JOINED.

`anti_join(x, y, by = NULL, ...)`
 Return rows of x that do not have a match in y. USEFUL TO SEE WHAT WILL NOT BE JOINED.

summarize() examples

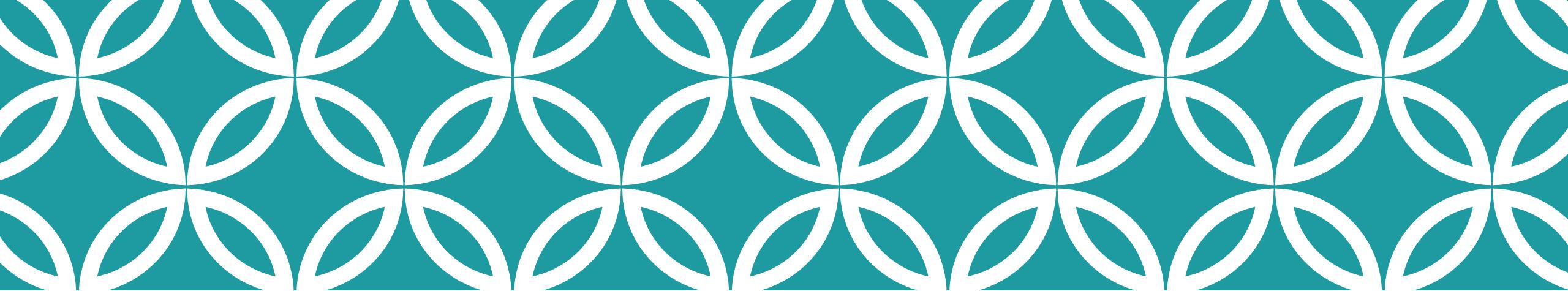
- Last pandemic day (in data)
- Median turnaround time



Your Turn 2

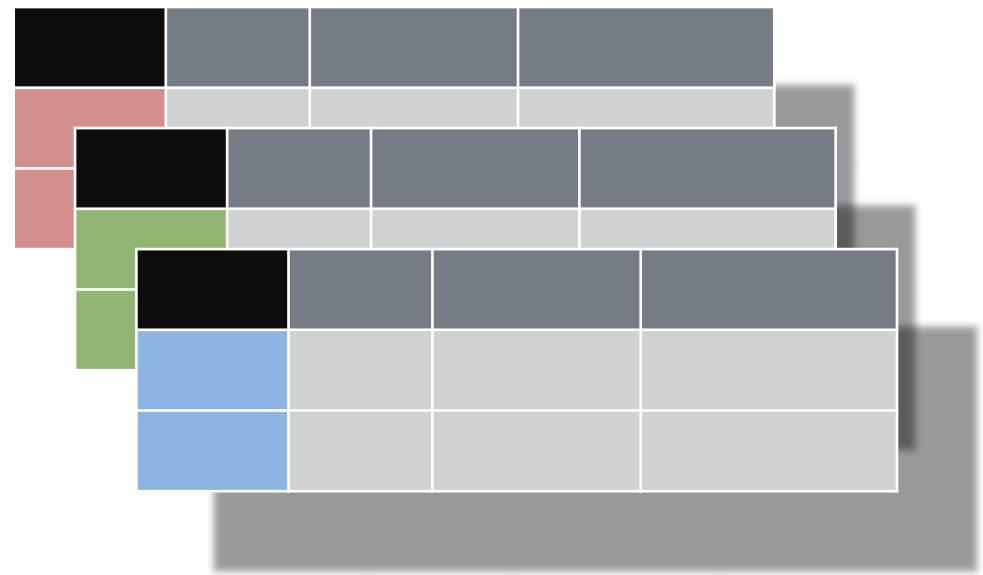
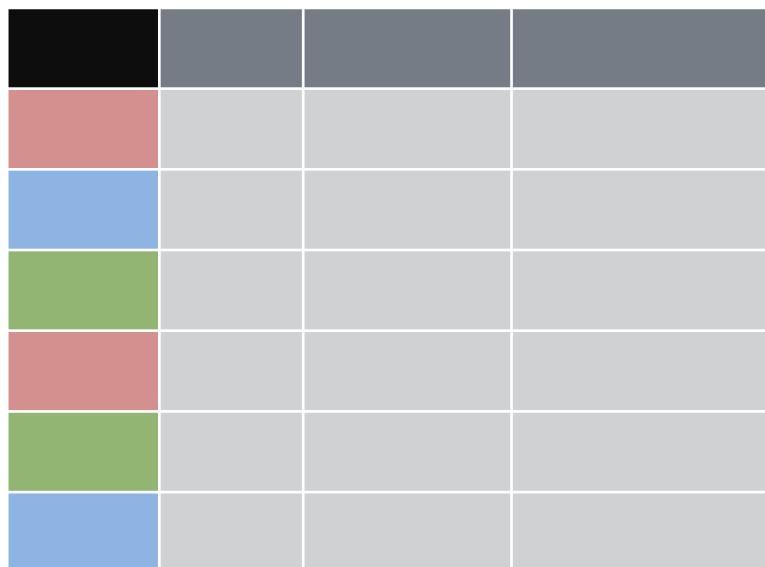
Consider:

How would you calculate the median number of orders per day?



group_by()

group_by()



group_by()

- *Grouping observations based on a specific variable's values*

```
covid_testing %>%  
  group_by(variable)
```

name of variable
to group by



group_by()

- *Group observations by pan_day*

```
covid_testing %>%  
  group_by(pan_day)
```

```
# A tibble: 15,524 x 17  
# Groups:   pan_day [102]  
  mrn first_name last_name gender pan_day  
  <dbl> <chr>     <chr>    <chr>   <dbl>  
1 5.00e6 jhezane   westerli... female      4  
2 5.00e6 penny     targaryen female      7  
3 5.01e6 grunt     rivers    male       7  
4 5.01e6 melisandre swyft    female      8  
5 5.01e6 rolley    karstark male       8
```

group_by()

- *Group observations by `pan_day` and `clinic_name`*

```
covid_testing %>%  
  select(mrn, pan_day, clinic_name) %>%  
  group_by(pan_day, clinic_name)
```

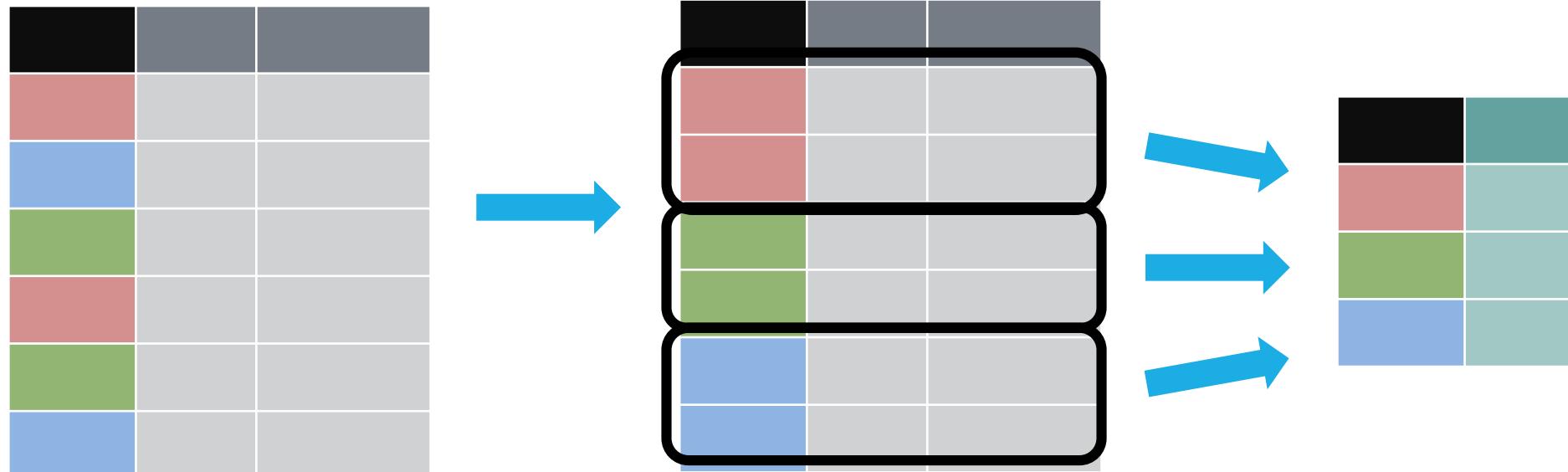
```
# A tibble: 15,524 x 3  
# Groups:   pan_day, clinic_name [2,526]  
  mrn pan_day clinic_name  
  <dbl> <dbl> <chr>  
1 5001412     4 inpatient ward a  
2 5000533     7 clinical lab  
3 5009134     7 clinical lab  
4 5008518     8 clinical lab  
5 5008967     8 emergency dept
```



```
group_by() %>% summarize()
```

`|group_by() %>% summarize()`

Make summaries of your data *by group*



group_by() %>% summarize()

- Make summaries of your data

```
covid_testing %>%  
  summarize(order_count = n())
```

| mrn | pan_day | order_count |
|---------|---------|-------------|
| 5001412 | 4 | 15524 |
| 5000533 | 7 | |
| 5009134 | 7 | |
| 5008518 | 8 | |



group_by() %>% summarize()

- Make summaries of your data

```
covid_testing %>%  
  group_by(pan_day) %>%  
  summarize(order_count = n())
```

| mrn | pan_day |
|---------|---------|
| 5001412 | 4 |
| 5000533 | 7 |
| 5009134 | 7 |
| 5008518 | 8 |



| pan_day | order_count |
|---------|-------------|
| 4 | 1 |
| 7 | 2 |
| 8 | 3 |
| 9 | 9 |

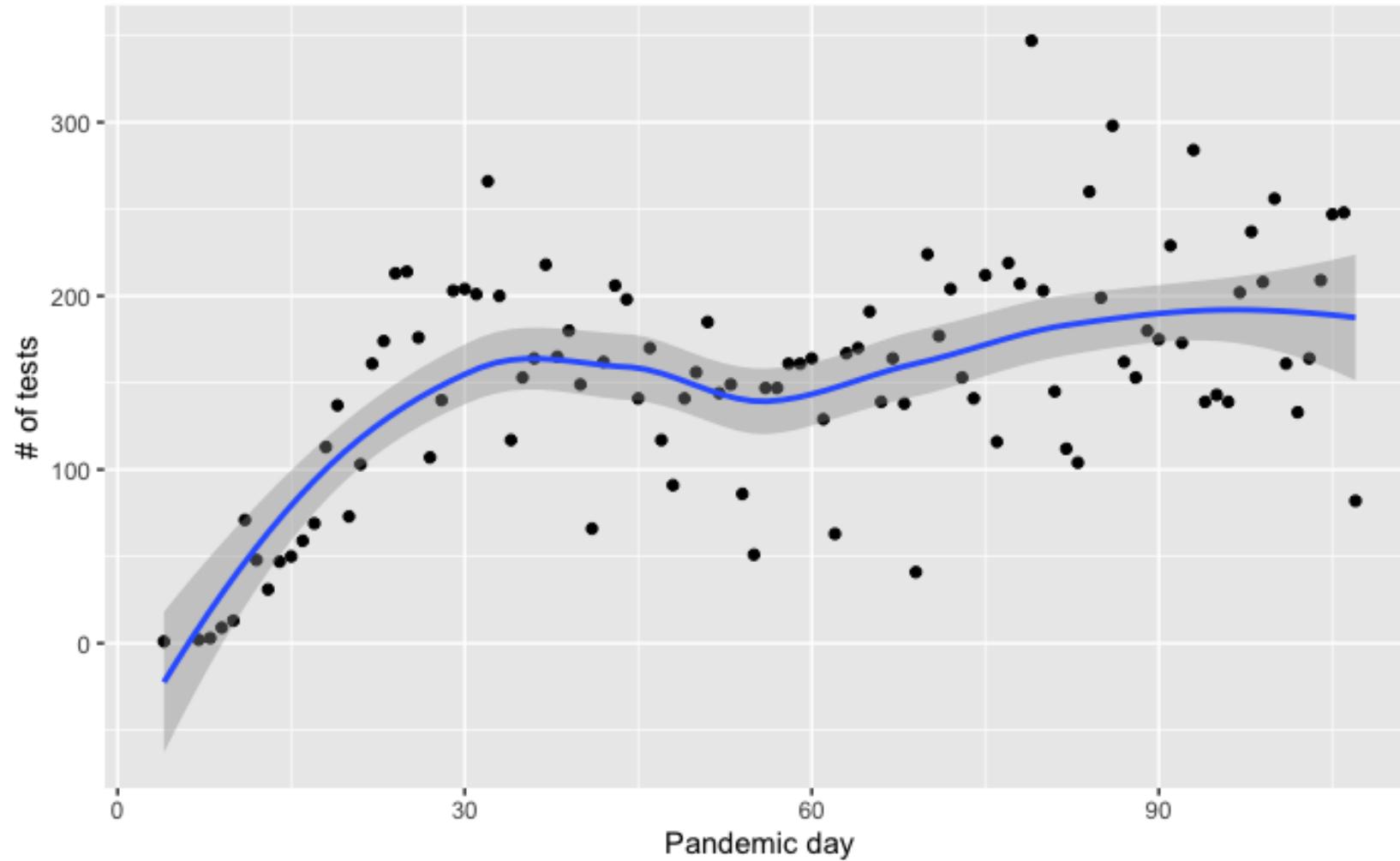


Your Turn 3

Calculate:

- a) The median turnaround time for each day
- b) (**Extra**) The median number of orders per day

group_by() %>% summarize(): Example



Stats: Tests for association

Q: Is there an association between insurance and SARS-CoV-2 RT-PCR positivity?

| payor_group_fac
<chr> | negative
<int> | positive
<int> |
|--------------------------|-------------------|-------------------|
| commercial | 3549 | 86 |
| government | 3318 | 242 |
| other | 309 | 17 |
| unassigned | 7182 | 520 |

4 rows



```
data %>%  
  fisher.test(simulate.p.value = T)
```

Data wrangling - 1

function that flexibly
assigns values

```
covid_testing_2 <- covid_testing %>%  
  mutate(payor_group_fac = case_when(  
    is.na(payor_group) ~ "unassigned",  
    payor_group %in% c("charity care",  
      "medical assistance",  
      "self pay",  
      "other") ~ "other",  
    TRUE ~ payor_group))  
  ) %>%  
  filter(result %in% c("positive", "negative"))
```



Data wrangling - 2

```
# Generate counts
tmp_table_tall <- covid_testing_2 %>%
  group_by(payor_group_fac, result) %>%
  summarize(n = n()) %>%
  ungroup()
tmp_table_tall

# Pivot from tall to wide table
tmp_table_wide <- tmp_table_tall %>%
  spread(key = "result", value = "n")
tmp_table_wide
```

Remove groupings

Maps key values to separate columns



Testing for association

| payor_group_fac
<chr> | negative
<int> | positive
<int> |
|--------------------------|-------------------|-------------------|
| commercial | 3549 | 86 |
| government | 3318 | 242 |
| other | 309 | 17 |
| unassigned | 7182 | 520 |

4 rows



```
data %>%  
  fisher.test(simulate.p.value = T)
```



Fisher's Exact Test for Count Data with simulated p-value (based on 2000 replicates)

```
data: .  
p-value = 0.0004998  
alternative hypothesis: two.sided
```



Testing for association: Example

Use `fisher.test()` to estimate the relative odds of a positive test result for patients with government insurance compared to commercial insurance?





What Else?

Logistic regression

```
tmp_fit <- glm(result_fac ~ payor_group_fac + age, # model formula  
                 data = tmp, # dataset  
                 family = "binomial") # type of model  
  
summary(tmp_fit)
```

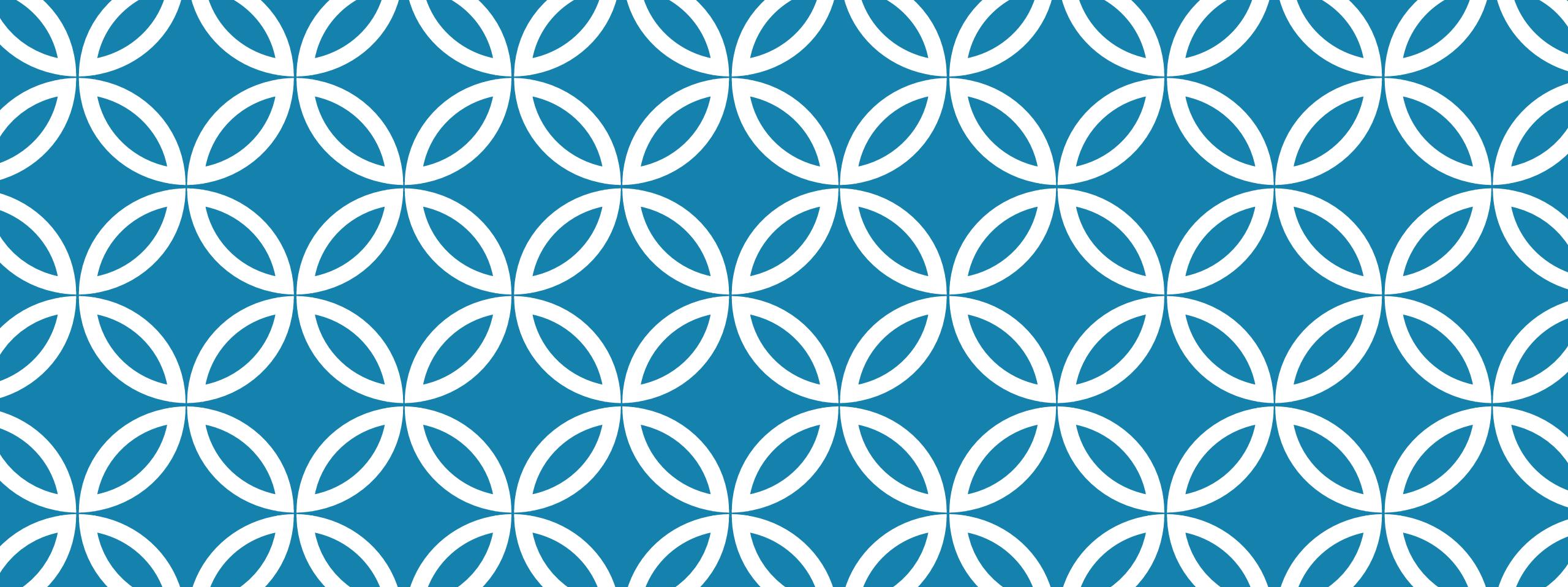


Goals

1. Learn how to summarize data
2. Appreciate how to perform statistical analyses

Objectives

1. Calculate a summary statistic for a variable using `summarize`
2. Calculate of summary statistic for a variable separately for a group of observations, using `group_by` and `summarize`
3. Describe the process of performing a simple test for association



Advanced Reporting

Session 7
Patrick Mathias

| July 16 2020 | Session | Instructor |
|-------------------|--|-------------------|
| 1:00 pm - 1:30 pm | Instructor Introductions, Introduction to technology | Amrom Obstfeld |
| 1:30 pm - 2:15 pm | Introduction to R and RStudio | Joe Rudolf |
| 2:30 pm - 3:15 pm | Reproducible Reporting | Patrick Mathias |
| 3:30 pm - 5:00 pm | Data Visualization | Stephan Kadauke |
| July 17 2020 | | |
| 1:00 pm - 2:30 pm | Data Transformation | Amrom Obstfeld |
| 2:45 pm - 4:15 pm | Statistical Analysis | Dan Herman |
| 4:30 pm - 5:00 pm | Advanced Reporting | Patrick Mathias |

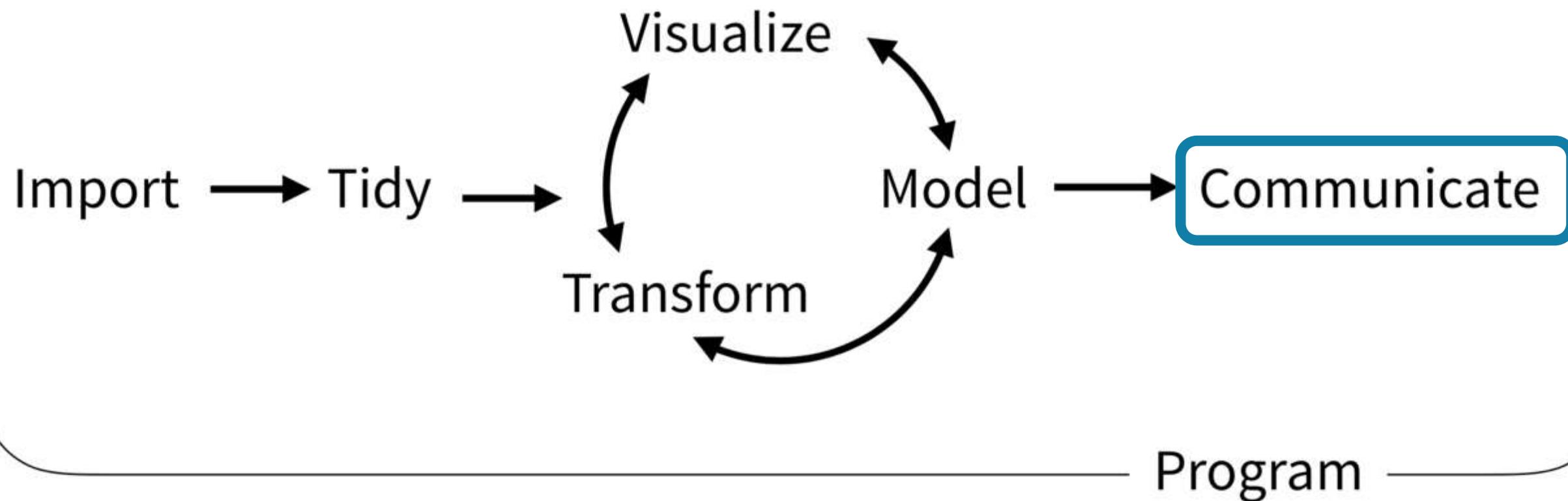
Goals

1. Build R Markdown reports using formatting outputs beyond standard document formats

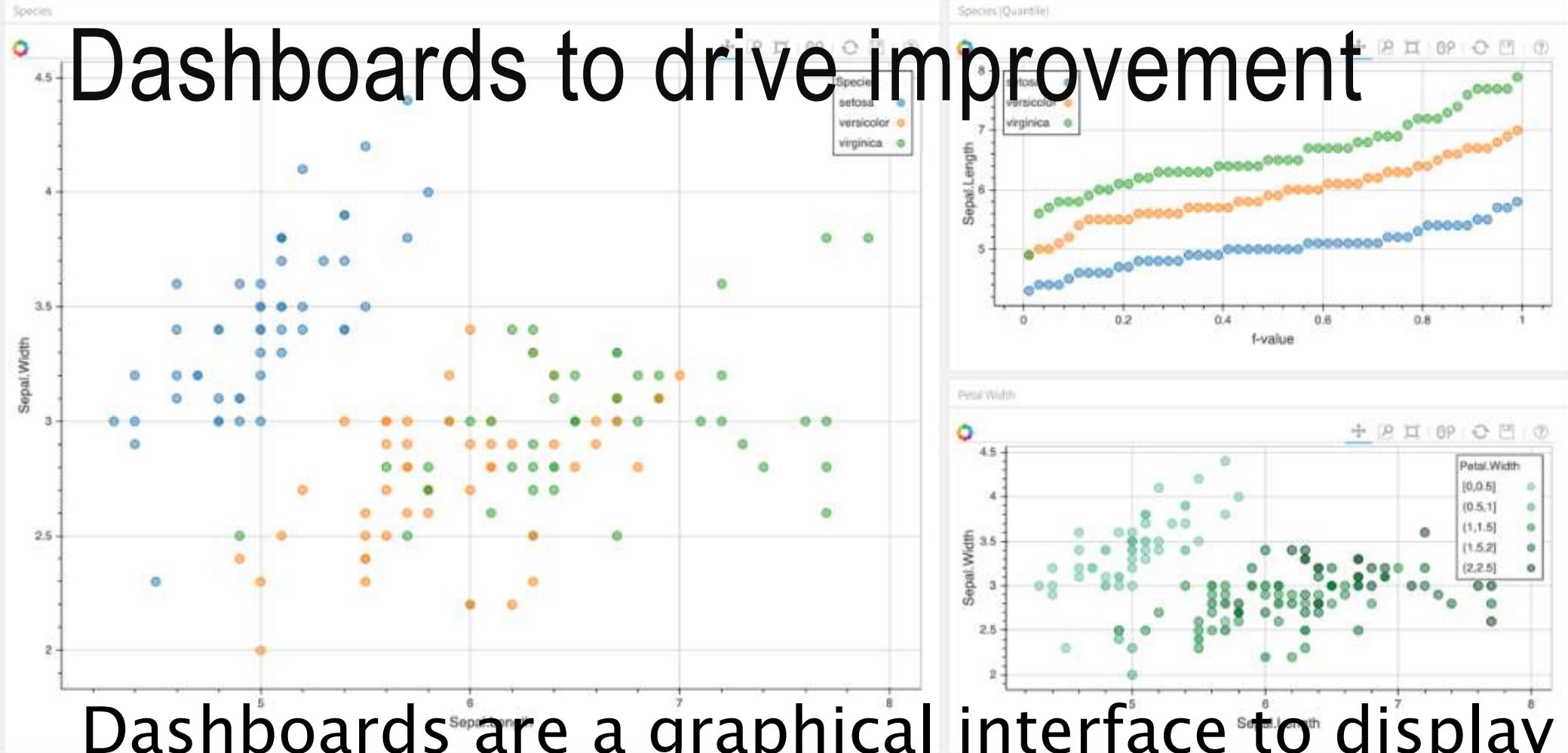
Objectives

1. Format a flexdashboard to improve display of multiple plots
2. Convert a static plot into an interactive plot

Typical Data Science Pipeline



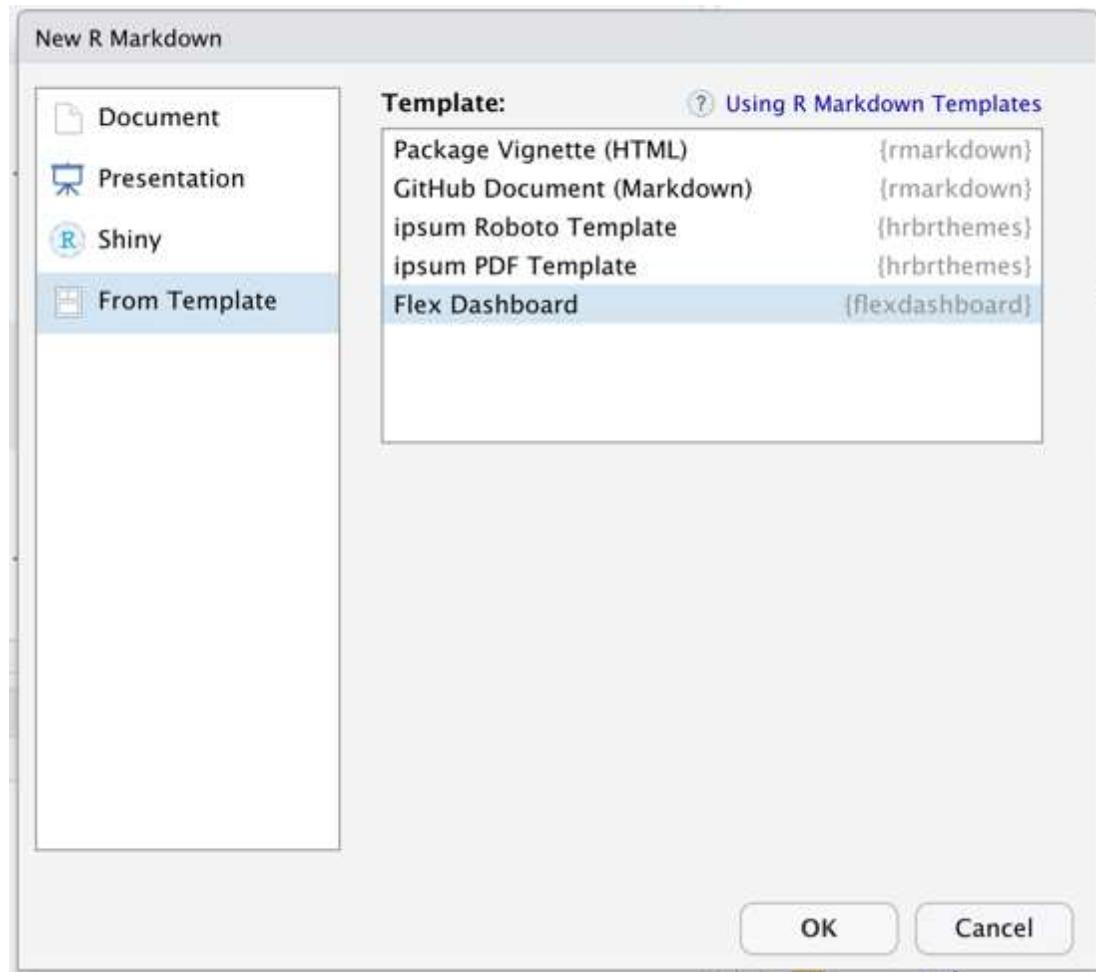
From R Markdown to Quick and Painless Dashboards



Dashboards are a graphical interface to display key performance indicators or other metrics

Intended to represent multiple pieces of information at a glance

flexdashboard provides easy dashboard templates for reporting



Produces HTML file that can be opened on web browsers

Or deployed on existing web server

Provides row or column based layouts

Get started by running:
`install.packages("flexdashboard")`

<https://rmarkdown.rstudio.com/flexdashboard/>

Untitled1 x ABC Knit Insert Run R Markdown

```
1 ---  
2 title: "Untitled"  
3 output:  
4   flexdashboard::flex_dashboard: ← flexdashboard output  
5     orientation: columns ← layout page by columns  
6     vertical_layout: fill  
7 ---  
8  
9 `r setup, include=FALSE}  
10 library(flexdashboard)  
11 ...  
12  
13 Column {data-width=650} define width  
14 -----  
15  
16 `## Chart A` title for chart ← delimits separate  
17 columns  
18 `r`  
19 ...  
20  
21  
22 Column {data-width=350}  
23 -----  
24  
25 ...
```

```
1 ---  
2 title: "Column Orientation"  
3 output: flexdashboard::flex_dashboard  
4 ---  
5  
6 Column  
7 -----  
8 |  
9 ### Chart 1  
10 ````{r}  
11 ...  
12  
13 Column  
14 -----  
15  
16  
17 ### Chart 2  
18 ````{r}  
19 ...  
20  
21  
22 ### Chart 3  
23  
24 ````{r}  
25 ...  
26
```

Chart 1

Chart 2

Chart 3

```
1 ---  
2 title: "Row Orientation"  
3 output:  
4   flexdashboard::flex_dashboard:  
5     orientation: rows  
6 ---  
7  
8 Row  
9 -----  
10  
11 ### Chart 1  
12  
13 `r`  
14  
15  
16 Row  
17 -----  
18  
19 ### Chart 2  
20  
21 `r`  
22  
23  
24 ### Chart 3  
25  
26 `r`  
27  
28
```

Chart 1

Chart 2

Chart 3

```
1 ---  
2 title: "Chart Stack (Scrolling)"  
3 output:  
4   flexdashboard::flex_dashboard:  
5     vertical_layout: scroll  
6 ---  
7  
8 ### Chart 1  
9  
10 `r`  
11 ...  
12  
13 ### Chart 2  
14  
15 `r`  
16 ...  
17  
18 ### Chart 3  
19  
20 `r`  
21 ...  
22  
23  
24  
25
```

Chart 1

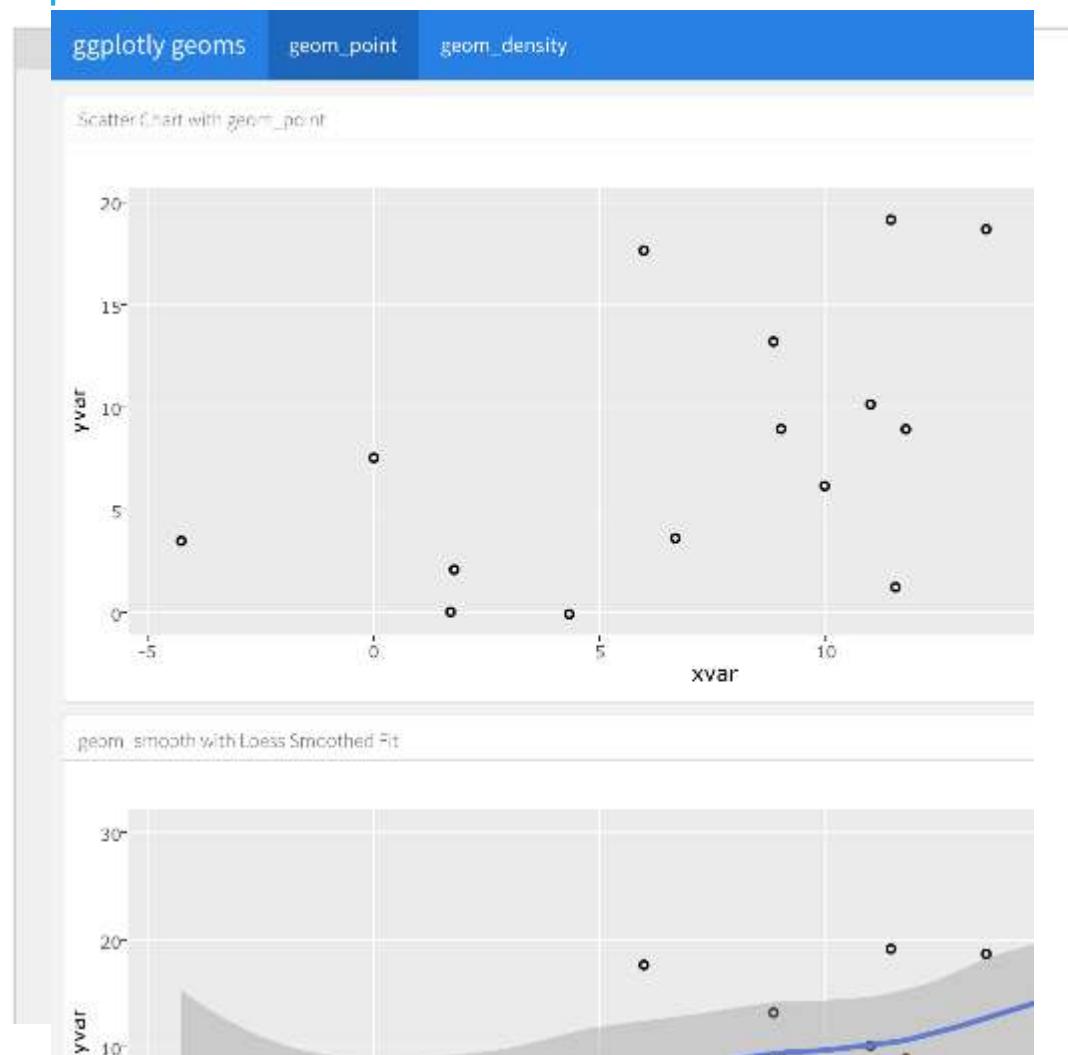
Chart 2

Chart 3

Your Turn #1

1. Open “07 – Advanced Reporting.Rmd” to work with a draft COVID-19 flexdashboard and run the setup chunk. Knit the document to see the dashboard output.
2. The “Test Volumes Over Time” plot could show additional information regarding positive tests. Add fill to your barplot to show the result field in addition to overall test volume by day. Run that code chunk to see the output.
3. Too much information is crunched on the right side. Change the layout from columns to a row orientation. The 2nd and 3rd plots (Turnaround Times and Cycle Thresholds) should appear on the 2nd row.

Customization



Bootswatch Themes ▾ Download ▾ Help Blog

Cerulean

A calm blue sky

Primary Secondary Success Info Warning

PREVIEW DOWNLOAD

Bootswatch Themes ▾ Download ▾ Help Blog

Cosmo

An ode to Metro

Primary Secondary Success Info Warning Danger

PREVIEW DOWNLOAD

Bootswatch Themes ▾ Download ▾ Help Blog

Darkly

Flatly in night mode

Primary Secondary Success Info Warning

PREVIEW DOWNLOAD

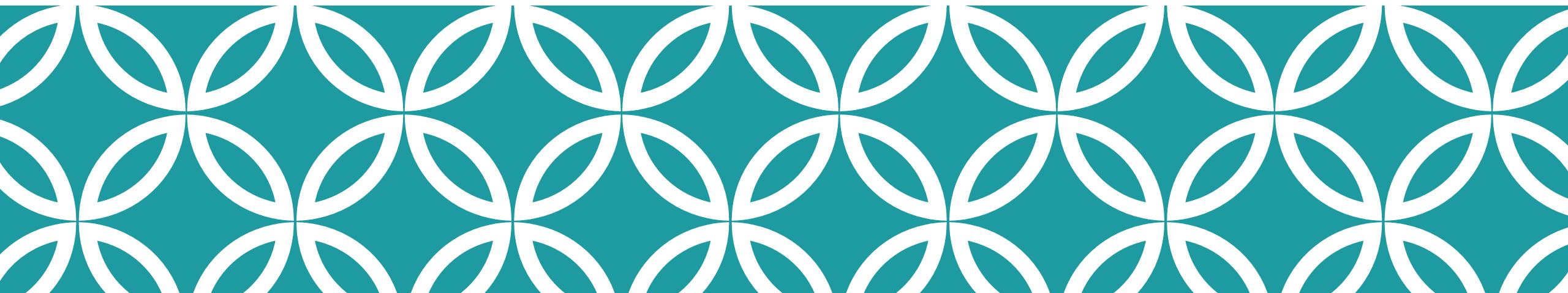
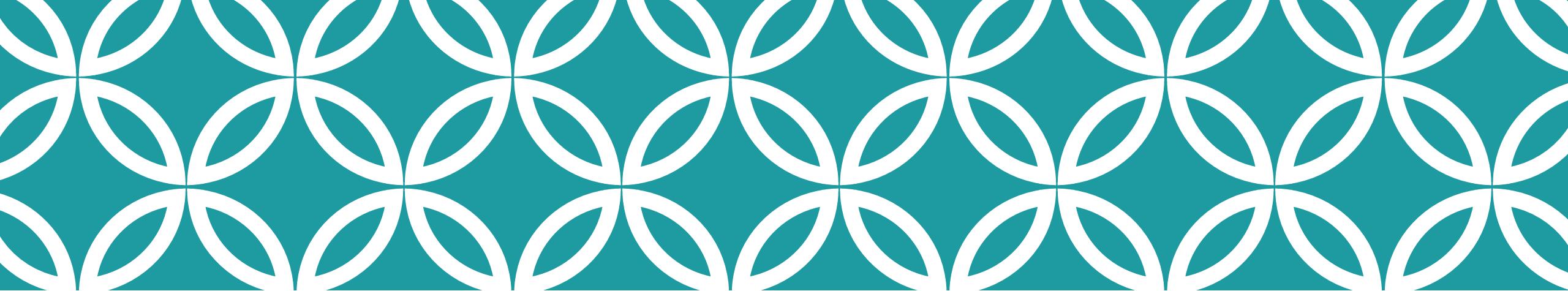
Bootswatch Themes ▾ Download ▾ Help Blog

Flatly

Flat and modern

Primary Secondary Success Info Warning

PREVIEW DOWNLOAD



Making plots interactive

htmlwidgets for R support interactive visuals

Packages using htmlwidgets use R code to call Javascript visualization libraries
(<http://www.htmlwidgets.org/>)

Use one line of code to convert a static plot into an interactive one

Plotly package converts ggplot with a simple command

To use Plotly install the plotly package using the following command:

```
install.packages("plotly")
```

Examples of visualizations at Plotly website:

<https://plotly.com/r/>

Store plot as object and add one line to make interactive

```
plot_name <- ggplot(data = data_frame) +  
  geom_function(mapping = aes(mappings))  
ggplotly(plot_name)
```

Your Turn #2

1. Load the `plotly` package in your setup chunk
2. Convert each of the plots into an interactive plot by storing the `ggplot` in an object and using the `ggplotly()` function.
3. Knit the dashboard and hover over the interactive plots.

Other options for interactive plots

Other interactive plot packages:

- rbokeh
- Highcharter

Time series graphs with dygraphs package

Maps with leaflet package

Interactive tables with one line

DataTables library quickly converts tables into interactive element

DT package in R

Use datatable() function on a data frame

- Filter number of entries
- Search entries
- Sort by column

datatable example

```
datatable(head(iris), class = 'cell-border stripe')
```

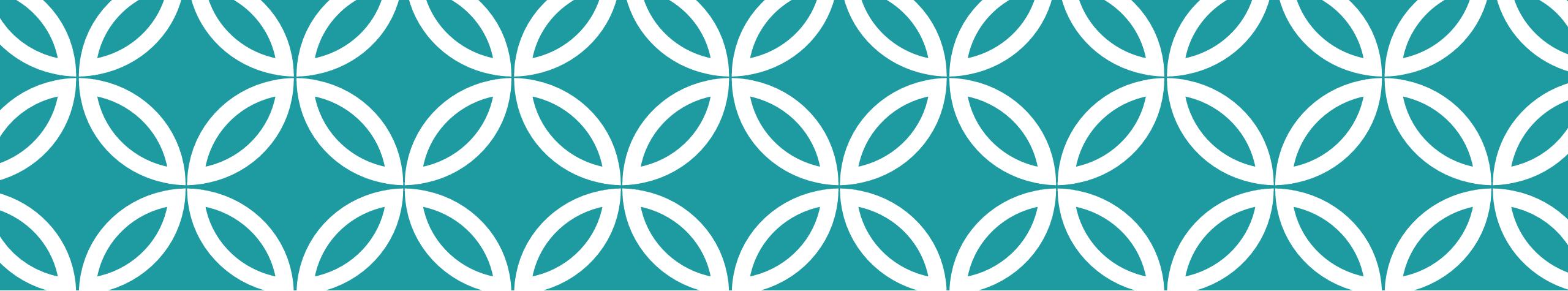
Show entries

Search:

| | Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species |
|---|--------------|-------------|--------------|-------------|---------|
| 1 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 2 | 4.9 | 3 | 1.4 | 0.2 | setosa |
| 3 | 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 4 | 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 5 | 5 | 3.6 | 1.4 | 0.2 | setosa |
| 6 | 5.4 | 3.9 | 1.7 | 0.4 | setosa |

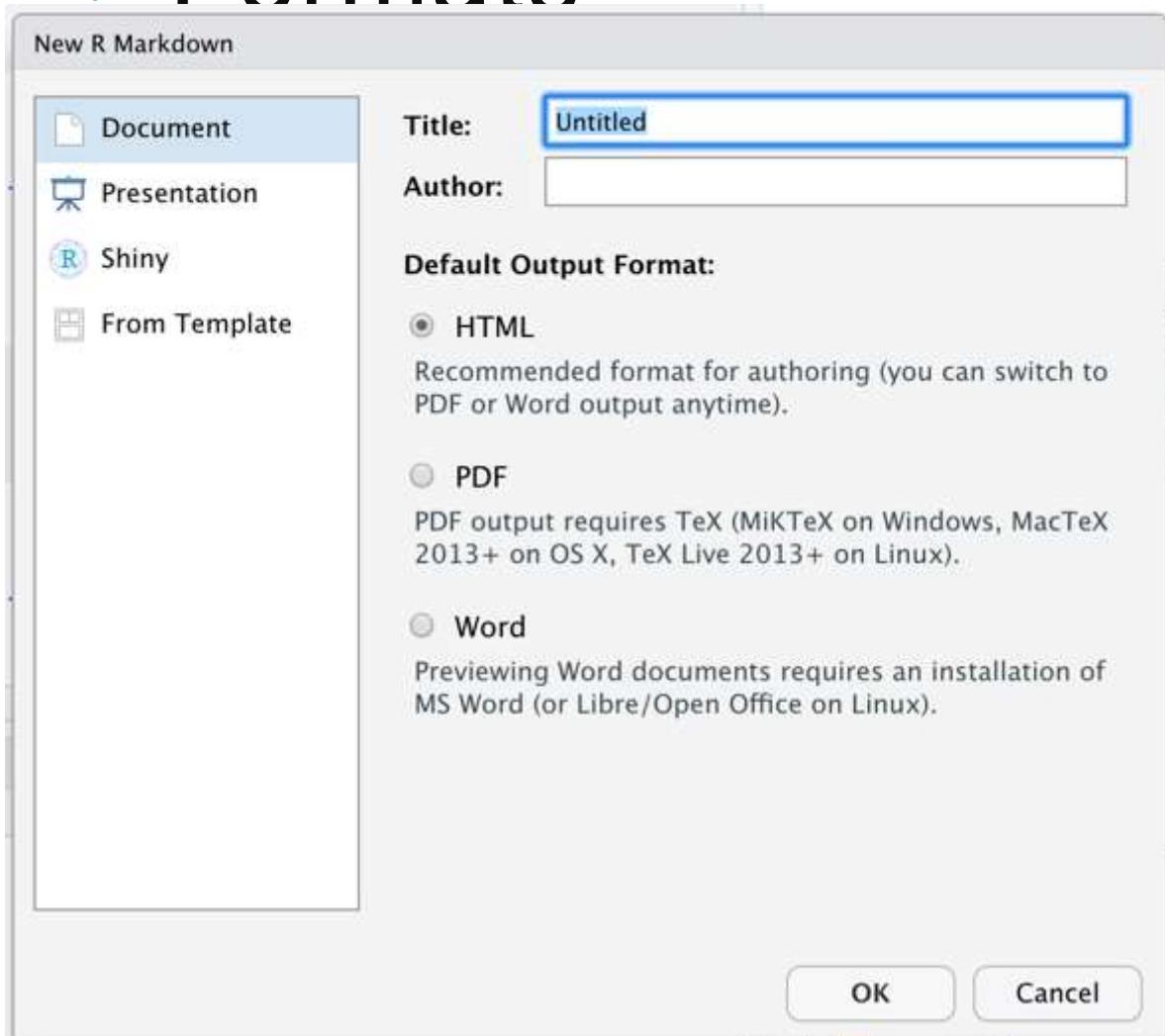
Showing 1 to 6 of 6 entries

Previous Next



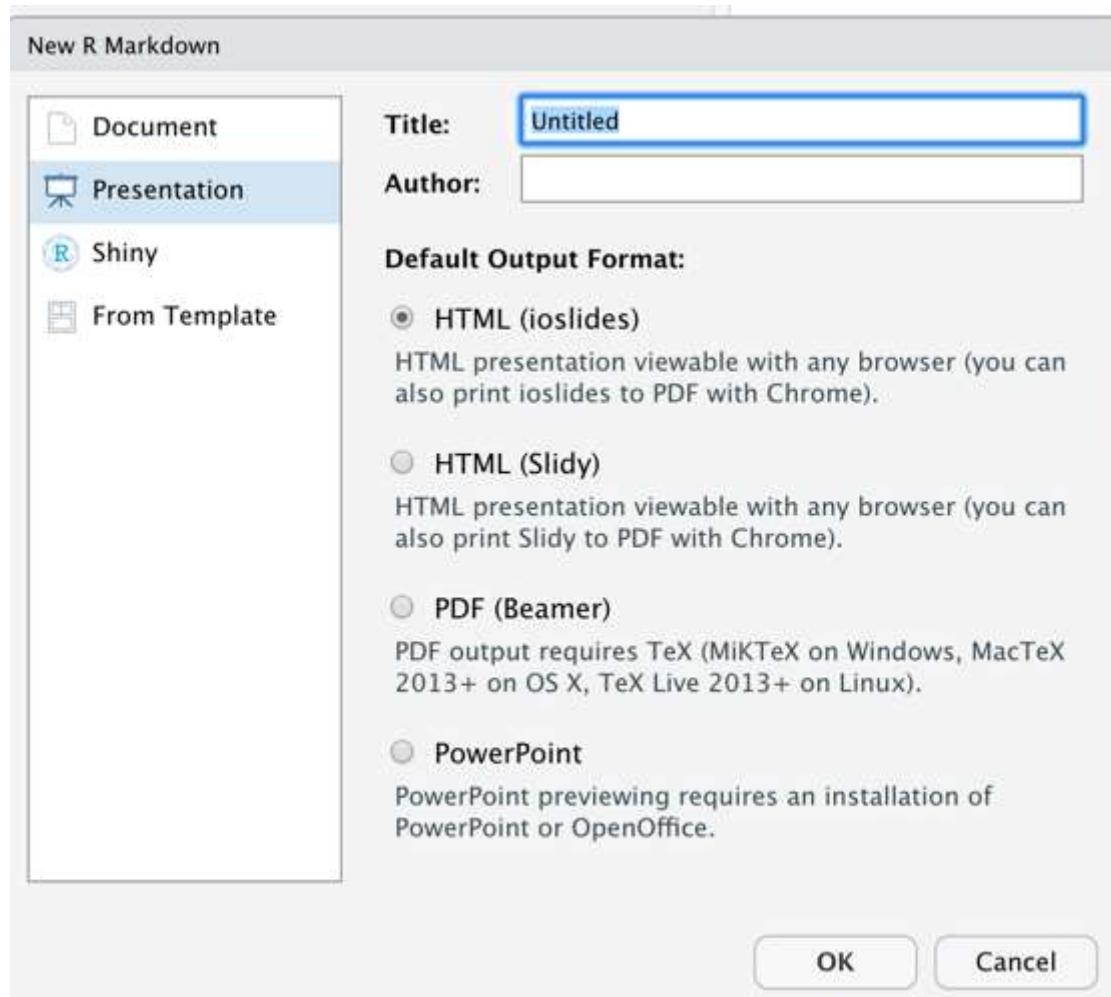
What Else?

Standard Markdown Reporting Formats



- HTML file – open with any web browser
- PDF – requires LaTeX dependencies
 - `install.packages('tinytex')`
 - `tinytex::install_tinytex()`
- Word – default format for collaborating with those who aren't familiar with R

Formats to go straight from code to slides



Multiple HTML formats create webpage that's advanceable like slides

PDF presentation uses LaTeX in the background

Powerpoint produces simple slides

The screenshot shows an RStudio interface with an R Markdown file open. The code editor contains the following R Markdown code:

```
1 - ---
2 title: "Untitled"
3 output: powerpoint_presentation
4 ---
5
6 ```{r setup, include=FALSE}
7 knitr::opts_chunk$set(echo = FALSE)
8
9
10 ## R Markdown
11
12 This is an R Markdown presentation. 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>.
13
14 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.
15
16 ## Slide with Bullets
17
18 - Bullet 1
19 - Bullet 2
20 - Bullet 3
21
22 ## Slide with R Output
23
24 ```{r cars, echo = TRUE}
25 summary(cars)
26
27
28 ## Slide with Plot
29
30 ```{r pressure}
31 plot(pressure)
32
33
34
```

Annotations with blue arrows point to specific parts of the code:

- An arrow points to the line `output: powerpoint_presentation` with the text "Output format".
- An arrow points to the line `## Slide with Bullets` with the text "Each slide has its own header".

R Markdown

This is an R Markdown presentation. 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>.

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.

Example output slide

Slide with Bullets

- Bullet 1
- Bullet 2
- Bullet 3

Example output slide

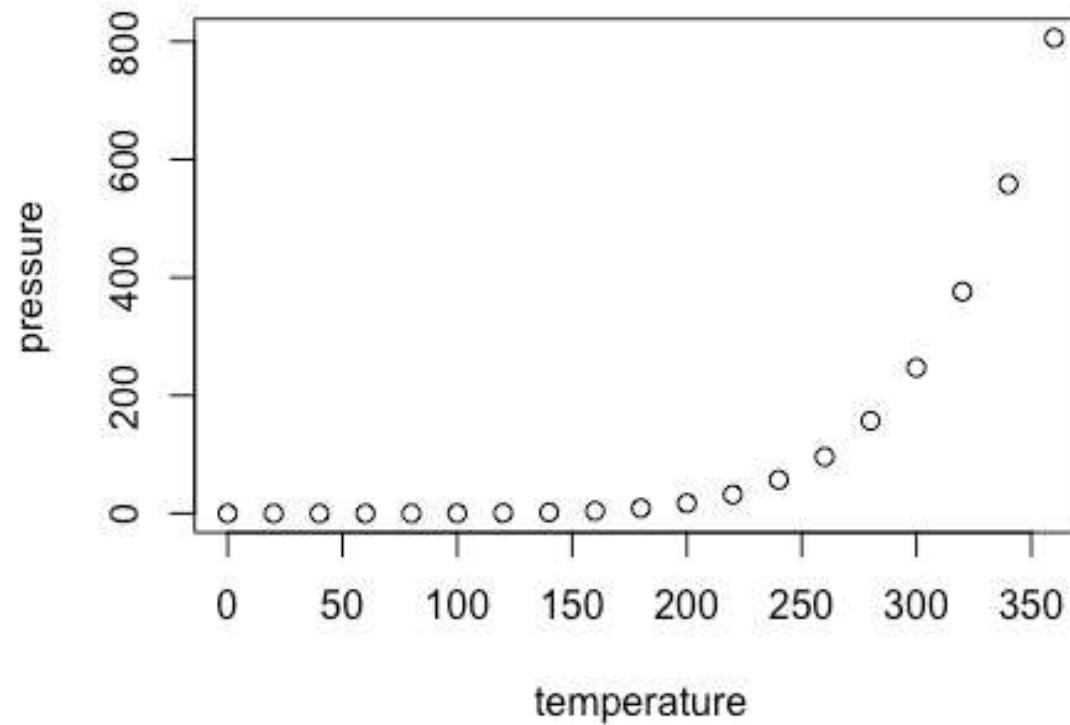
Slide with R Output

```
summary(cars)
```

| | speed | dist |
|------------|-------|----------------|
| ## Min. | : 4.0 | Min. : 2.00 |
| ## 1st Qu. | :12.0 | 1st Qu.: 26.00 |
| ## Median | :15.0 | Median : 36.00 |
| ## Mean | :15.4 | Mean : 42.98 |
| ## 3rd Qu. | :19.0 | 3rd Qu.: 56.00 |
| ## Max. | :25.0 | Max. :120.00 |

Example output slide

Slide with Plot



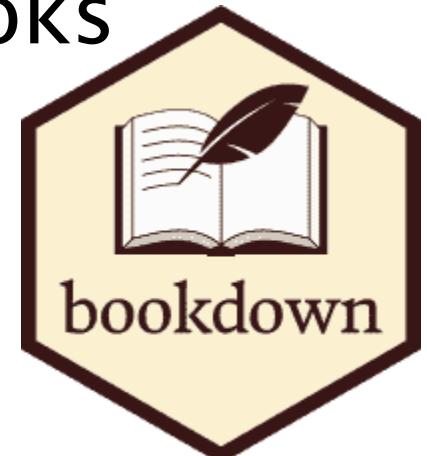
Example output slide

Books and longer documents also generated from R Markdown

Can generate printer ready books and ebooks

Supports LaTeX features such as equations

Generates blog formatted websites



<https://github.com/rstudio/bookdown>

<https://bookdown.org/yihui/bookdown/>

<https://bookdown.org/yihui/blogdown/>

Goals

1. Build R Markdown reports using formatting outputs beyond standard document formats

Objectives

1. Format a flexdashboard to improve display of multiple plots
2. Convert a static plot into an interactive plot

Appendix

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

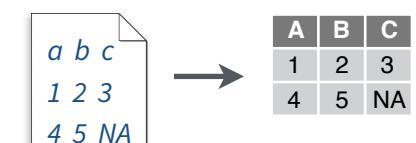
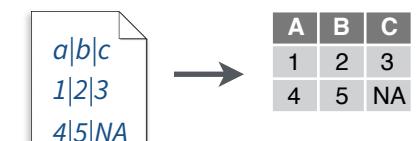
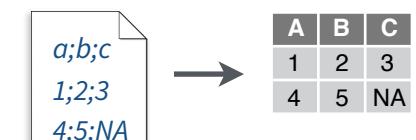
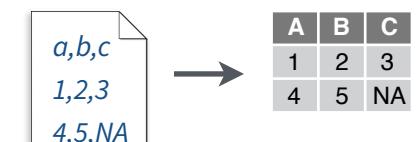
```
write_tsv(x, path, na = "NA", append = FALSE,  
          col_names = !append)
```



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())
```



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))
```

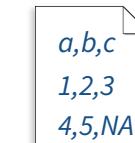
```
write_file(x = "a b c\n1 2 3\n4 5 NA", path = "file.fwf")
```

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



Example file

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

| | | |
|---|---|----|
| 1 | 2 | 3 |
| 4 | 5 | NA |

Skip lines

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

| | | |
|---|---|----|
| A | B | C |
| 1 | 2 | 3 |
| 4 | 5 | NA |

No header

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

| | | |
|---|---|----|
| A | B | C |
| 1 | 2 | 3 |
| 4 | 5 | NA |

Read in a subset

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

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

Provide header

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

| | | |
|----|---|----|
| 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(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)`

Tibbles - an enhanced data frame

The **tibble** package provides a new S3 class for storing tabular data, the **tibble**. Tibbles inherit the data frame class, but improve three behaviors:

- **Subsetting** - [always returns a new tibble, [[and \$ always return a vector.
- **No partial matching** - You must use full column names when subsetting
- **Display** - When you print a tibble, R provides a concise view of the data that fits on one screen

| # A tibble: 234 × 6 | manufacturer | model | displ | cyl | trans |
|--|--------------|-------|-------|----------|-------|
| 1 audi | a4 | 1.80 | 4 | auto(l4) | |
| 2 audi | a4 | 1.80 | 4 | auto(l4) | |
| 3 audi | a4 | 2.00 | 4 | auto(l4) | |
| 4 audi | a4 | 2.00 | 4 | auto(l4) | |
| 5 audi | a4 | 2.00 | 4 | auto(l4) | |
| 6 audi | a4 | 2.00 | 4 | auto(l4) | |
| 7 audi | a4 | 3.10 | 6 | quattro | |
| 8 audi | a4 | 3.10 | 6 | quattro | |
| 9 audi | a4 | 3.10 | 6 | quattro | |
| 10 audi | a4 | 3.10 | 6 | quattro | |
| ... with 224 more rows, and 3 more variables: year <int>, cyl <int>, trans <chr> | | | | | |

tibble display

| country | 1999 | 2000 |
|---------|------|------|
| A | 0.7K | 2K |
| B | 37K | 80K |
| C | 212K | 213K |

A large table to display

data frame display

- Control the default appearance with options:
`options(tibble.print_max = n, tibble.print_min = m, tibble.width = Inf)`
- View full data set with **View()** or **glimpse()**
- Revert to data frame with **as.data.frame()**

CONSTRUCT A TIBBLE IN TWO WAYS

| | | |
|---------------------|---|---|
| tibble(...) | Construct by columns.
<code>tibble(x = 1:3, y = c("a", "b", "c"))</code> | Both make this tibble |
| tribble(...) | Construct by rows.
<code>tribble(~x, ~y, 1, "a", 2, "b", 3, "c")</code> | A tibble: 3 × 2
x y
1 a
2 b
3 c |

as_tibble(x, ...) Convert data frame to tibble.
enframe(x, name = "name", value = "value") Convert named vector to a tibble
is_tibble(x) Test whether x is a tibble.

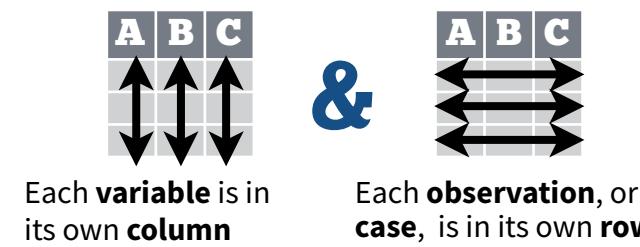


R Studio

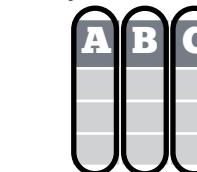
Tidy Data with tidyverse

Tidy data is a way to organize tabular data. It provides a consistent data structure across packages.

A table is tidy if:

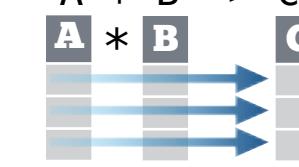


Tidy data:



Makes variables easy to access as vectors

$A * B \rightarrow C$



Preserves cases during vectorized operations

Reshape Data - change the layout of values in a table

Use **gather()** and **spread()** to reorganize the values of a table into a new layout.

gather(data, key, value, ..., na.rm = FALSE, convert = FALSE, factor_key = FALSE)

gather() moves column names into a **key** column, gathering the column values into a single **value** column.

table4a

| country | 1999 | 2000 |
|---------|------|------|
| A | 0.7K | 2K |
| B | 37K | 80K |
| C | 212K | 213K |



| country | year | cases |
|---------|------|-------|
| A | 1999 | 0.7K |
| B | 1999 | 37K |
| C | 1999 | 212K |
| A | 2000 | 2K |
| B | 2000 | 80K |
| C | 2000 | 213K |

key value

`gather(table4a, `1999`, `2000`, key = "year", value = "cases")`

table2

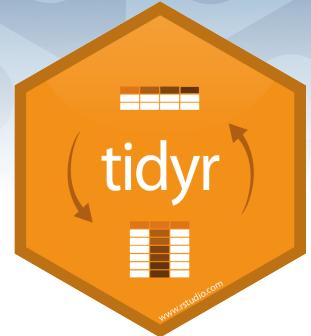
| country | year | type | count |
|---------|------|-------|-------|
| A | 1999 | cases | 0.7K |
| A | 1999 | pop | 19M |
| A | 2000 | cases | 2K |
| A | 2000 | pop | 20M |
| B | 1999 | cases | 37K |
| B | 1999 | pop | 172M |
| B | 2000 | cases | 80K |
| B | 2000 | pop | 174M |
| C | 1999 | cases | 212K |
| C | 1999 | pop | 1T |
| C | 2000 | cases | 213K |
| C | 2000 | pop | 1T |

key value

`spread(table2, type, count)`

Split Cells

Use these functions to split or combine cells into individual, isolated values.



separate(data, col, into, sep = "[^[:alnum:]]+", remove = TRUE, convert = FALSE, extra = "warn", fill = "warn", ...)

Separate each cell in a column to make several columns.

| country | year | rate | cases | pop |
|---------|------|----------|-------|-----|
| A | 1999 | 0.7K/19M | | |
| A | 2000 | 2K/20M | | |
| B | 1999 | 37K/172M | | |
| B | 2000 | 80K/174M | | |
| C | 1999 | 212K/1T | | |
| C | 2000 | 213K/1T | | |

`separate(table3, rate, sep = "/", into = c("cases", "pop"))`

separate_rows(data, ..., sep = "[^[:alnum:]].+", convert = FALSE)

Separate each cell in a column to make several rows.

| country | year | rate | cases | pop |
|---------|------|------|-------|-----|
| A | 1999 | 0.7K | 19M | |
| A | 1999 | 19M | | |
| A | 2000 | 2K | | |
| A | 2000 | 20M | | |
| B | 1999 | 37K | | |
| B | 1999 | 172M | | |
| B | 2000 | 80K | | |
| B | 2000 | 174M | | |
| C | 1999 | 212K | | |
| C | 1999 | 1T | | |
| C | 2000 | 213K | | |
| C | 2000 | 1T | | |

`separate_rows(table3, rate, sep = "/")`

unite(data, col, ..., sep = "_", remove = TRUE)

Collapse cells across several columns to make a single column.

| country | century | year |
|---------|---------|------|
| Afghan | 19 | 99 |
| Afghan | 20 | 00 |
| Brazil | 19 | 99 |
| Brazil | 20 | 00 |
| China | 19 | 99 |
| China | 20 | 00 |

`unite(table5, century, year, col = "year", sep = "")`

Handle Missing Values

drop_na(data, ...)

Drop rows containing NA's in ... columns.

| x1 | x2 |
|----|----|
| A | 1 |
| B | NA |
| C | NA |
| D | 3 |
| E | NA |

`drop_na(x, x2)`

fill(data, ..., .direction = c("down", "up"))

Fill in NA's in ... columns with most recent non-NA values.

| x1 | x2 |
|----|----|
| A | 1 |
| B | NA |
| C | NA |
| D | 3 |
| E | NA |

`fill(x, x2)`

replace_na(data, replace = list(), ...)

Replace NA's by column.

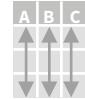
| x1 | x2 |
|----|----|
| A | 1 |
| B | NA |
| C | NA |
| D | 3 |
| E | NA |

`replace_na(x, list`

Data Transformation with dplyr :: CHEAT SHEET



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



Each **variable** is in its own **column**



Each **observation**, or **case**, is in its own **row**



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

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).

summary function

- `summarise(.data, ...)`
Compute table of summaries.
`summarise(mtcars, avg = mean(mpg))`
- `count(x, ..., wt = NULL, sort = FALSE)`
Count number of rows in each group defined by the variables in ... Also **tally()**.
`count(iris, Species)`

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.

- `mtcars %>% group_by(cyl) %>% summarise(avg = mean(mpg))`

`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.



`filter(.data, ...)` Extract rows that meet logical criteria. `filter(iris, Sepal.Length > 7)`



`distinct(.data, ..., .keep_all = FALSE)` Remove rows with duplicate values.
`distinct(iris, Species)`



`sample_frac(tbl, size = 1, replace = FALSE, weight = NULL, .env = parent.frame())` Randomly select fraction of rows.
`sample_frac(iris, 0.5, replace = TRUE)`



`slice(.data, ...)` Select rows by position.
`slice(iris, 10:15)`



`top_n(x, n, wt)` Select and order top n entries (by group if grouped data).
`top_n(iris, 5, Sepal.Width)`

Manipulate Variables

EXTRACT VARIABLES

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



`pull(.data, var = -1)` Extract column values as a vector. Choose by name or index.
`pull(iris, Sepal.Length)`



`select(.data, ...)` Extract columns as a table. Also **select_if()**.
`select(iris, Sepal.Length, Species)`

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

`mutate(.data, ...)`
Compute new column(s).
`mutate(mtcars, gpm = 1/mpg)`

`transmute(.data, ...)`
Compute new column(s), drop others.
`transmute(mtcars, gpm = 1/mpg)`

`mutate_all(.tbl, .funs, ...)` Apply funs to every column. Use with `funs()`. Also **mutate_if()**.
`mutate_all(faithful, funs(log(.), log2(.)))`
`mutate_if(iris, is.numeric, funs(log(.)))`

`mutate_at(.tbl, .cols, .funs, ...)` Apply funs to specific columns. Use with `funs()`, `vars()` and the helper functions for `select()`.
`mutate_at(iris, vars(-Species), funs(log(.)))`

`add_column(.data, ..., .before = NULL, .after = NULL)` Add new column(s). Also **add_count()**, **add_tally()**.
`add_column(mtcars, new = 1:32)`

`rename(.data, ...)` Rename columns.
`rename(iris, Length = Sepal.Length)`



Vector Functions

TO USE WITH MUTATE ()

mutate() and **transmute()** apply vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.

vectorized function

OFFSETS

dplyr::lag() - Offset elements by 1
dplyr::lead() - Offset elements by -1

CUMULATIVE AGGREGATES

dplyr::cumall() - Cumulative all()
dplyr::cumany() - Cumulative any()
cummax() - Cumulative max()
dplyr::cummean() - Cumulative mean()
cummin() - Cumulative min()
cumprod() - Cumulative prod()
cumsum() - Cumulative sum()

RANKINGS

dplyr::cume_dist() - Proportion of all values <= dplyr::dense_rank() - rank w ties = min, no gaps dplyr::min_rank() - rank with ties = min dplyr::ntile() - bins into n bins dplyr::percent_rank() - min_rank scaled to [0,1] dplyr::row_number() - rank with ties = "first"

MATH

+, -, *, /, ^, %/%, %% - arithmetic ops
log(), log2(), log10() - logs
<, <=, >, >=, !=, == - logical comparisons
dplyr::between() - x >= left & x <= right
dplyr::near() - safe == for floating point numbers

MISC

dplyr::case_when() - multi-case if_else()
iris %>% mutate(Species = case_when(
Species == "versicolor" ~ "versi",
Species == "virginica" ~ "virgi",
TRUE ~ Species))

dplyr::coalesce() - first non-NA values by element across a set of vectors
dplyr::if_else() - element-wise if() + else()
dplyr::na_if() - replace specific values with NA
pmax() - element-wise max()
pmin() - element-wise min()
dplyr::recode() - Vectorized switch()
dplyr::recode_factor() - Vectorized switch() for factors

Summary Functions

TO USE WITH SUMMARISE ()

summarise() applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

summary function

COUNTS

dplyr::n() - number of values/rows
dplyr::n_distinct() - # of uniques
sum(!is.na()) - # of non-NA's

LOCATION

mean() - mean, also mean(!is.na())
median() - median

LOGICALS

mean() - Proportion of TRUE's
sum() - # of TRUE's

POSITION/ORDER

dplyr::first() - first value
dplyr::last() - last value
dplyr::nth() - value in nth location of vector

RANK

quantile() - nth quantile
min() - minimum value
max() - maximum value

SPREAD

IQR() - Inter-Quartile Range
mad() - median absolute deviation
sd() - standard deviation
var() - variance

Row Names

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.

| A | B |
|---|---|
| 1 | a |
| 2 | b |
| 3 | c |

rownames_to_column()

Move row names into col.
a <- rownames_to_column(iris, var = "C")

| A | B | C |
|---|---|---|
| 1 | a | t |
| 2 | b | u |
| 3 | c | v |

column_to_rownames()

Move col in row names.
column_to_rownames(a, var = "C")

Also has_rownames(), remove_rownames()

Combine Tables

COMBINE VARIABLES

| | | | | | | |
|---|----------------------------------|---|---|----------------------------------|---|--|
| X | A B C
a t 1
b u 2
c v 3 | + | y | A B D
a t 3
b u 2
d w 1 | = | A B C A B D
a t 1 a t 3
b u 2 b u 2
c v 3 d w 1 |
|---|----------------------------------|---|---|----------------------------------|---|--|

Use **bind_cols()** to paste tables beside each other as they are.

bind_cols(...) Returns tables placed side by side as a single table.
BE SURE THAT ROWS ALIGN.

Use a "**Mutating Join**" to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

| | |
|---|--|
| A B C D
a t 1 3
b u 2 2
c v 3 NA | left_join(x, y, by = NULL,
copy = FALSE, suffix = c(".x", ".y"), ...) |
| | Join matching values from y to x. |

| | |
|---|---|
| A B C D
a t 1 3
b u 2 2
d w NA 1 | right_join(x, y, by = NULL, copy =
FALSE, suffix = c(".x", ".y"), ...) |
| | Join matching values from x to y. |

| | |
|-------------------------------|---|
| A B C D
a t 1 3
b u 2 2 | inner_join(x, y, by = NULL, copy =
FALSE, suffix = c(".x", ".y"), ...) |
| | Join data. Retain only rows with matches. |

| | |
|---|--|
| A B C D
a t 1 3
b u 2 2
c v 3 NA | full_join(x, y, by = NULL,
copy = FALSE, suffix = c(".x", ".y"), ...) |
| | Join data. Retain all values, all rows. |

Use **by = c("col1", "col2", ...)** to specify one or more common columns to match on.
left_join(x, y, by = "A")

Use a named vector, **by = c("col1" = "col2")**, to match on columns that have different names in each table.
left_join(x, y, by = c("C" = "D"))

Use **suffix** to specify the suffix to give to unmatched columns that have the same name in both tables.
left_join(x, y, by = c("C" = "D"), suffix = c("1", "2"))

COMBINE CASES

| | | | | |
|---|----------------------------------|---|---|-------------------------|
| X | A B C
a t 1
b u 2
c v 3 | + | y | A B C
C v 3
d w 4 |
|---|----------------------------------|---|---|-------------------------|

Use **bind_rows()** to paste tables below each other as they are.

| | |
|---|--|
| DF A B C
x a t 1
x b u 2
x c v 3
z c v 3
z d w 4 | bind_rows(..., .id = NULL) |
| | Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured) |

| | |
|----------------|-----------------------------------|
| A B C
c v 3 | intersect(x, y, ...) |
| | Rows that appear in both x and y. |

| | |
|-------------------------|----------------------------------|
| A B C
a t 1
b u 2 | setdiff(x, y, ...) |
| | Rows that appear in x but not y. |

| | |
|---|---|
| A B C
a t 1
b u 2
c v 3
d w 4 | union(x, y, ...) |
| | Rows that appear in x or y. (Duplicates removed). union_all() retains duplicates. |

Use **setequal()** to test whether two data sets contain the exact same rows (in any order).

EXTRACT ROWS

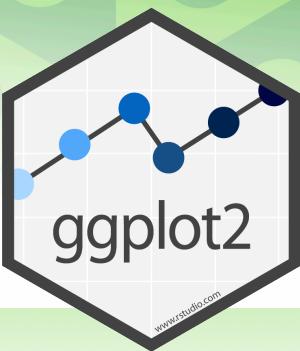
| | | | | | |
|---|----------------------------------|---|---|----------------------------------|---|
| X | A B C
a t 1
b u 2
c v 3 | + | y | A B D
a t 3
b u 2
d w 1 | = |
|---|----------------------------------|---|---|----------------------------------|---|

Use a "**Filtering Join**" to filter one table against the rows of another.

| | |
|-------------------------|---|
| A B C
a t 1
b u 2 | semi_join(x, y, by = NULL, ...) |
| | Return rows of x that have a match in y. USEFUL TO SEE WHAT WILL BE JOINED. |

| | |
|----------------|--|
| A B C
c v 3 | anti_join(x, y, by = NULL, ...) |
| | Return rows of x that do not have a match in y. USEFUL TO SEE WHAT WILL NOT BE JOINED. |

Data Visualization with ggplot2 :: CHEAT SHEET



Basics

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



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.



Complete the template below to build a graph.

```
ggplot (data = <DATA>) +
<GEOM_FUNCTION>(mapping = aes(<MAPPINGS>),
stat = <STAT>, position = <POSITION>) +
<COORDINATE_FUNCTION> +
<FACET_FUNCTION> +
<SCALE_FUNCTION> +
<THEME_FUNCTION>
```

↑ required
Not required, sensible defaults supplied

ggplot(data = mpg, **aes**(x = cty, y = hwy)) Begins a plot that you finish by adding layers to. Add one geom function per layer.

aesthetic mappings **data** **geom**

qplot(x = cty, y = hwy, data = mpg, geom = "point") Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

last_plot() Returns the last plot

ggsave("plot.png", **width** = 5, **height** = 5) Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.

Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

GRAPHICAL PRIMITIVES

- a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))
- a + geom_blank()**
(Useful for expanding limits)
- b + geom_curve(aes(yend = lat + 1, xend = long + 1), curvature = 1)** - x, yend, y, yend, alpha, angle, color, curvature, linetype, size
- a + geom_path(lineend = "butt", linejoin = "round", linemitre = 1)** - x, y, alpha, color, group, linetype, size
- a + geom_polygon(aes(group = group))** - x, y, alpha, color, fill, group, linetype, size
- b + geom_rect(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1))** - xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size
- a + geom_ribbon(aes(ymin = unemploy - 900, ymax = unemploy + 900))** - x, ymax, ymin, alpha, color, fill, group, linetype, size

LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

- b + geom_abline(aes(intercept = 0, slope = 1))**
- b + geom_hline(aes(yintercept = lat))**
- b + geom_vline(aes(xintercept = long))**
- b + geom_segment(aes(yend = lat + 1, xend = long + 1))**
- b + geom_spoke(aes(angle = 1:1155, radius = 1))**

ONE VARIABLE continuous

- c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
- c + geom_area(stat = "bin")** - x, y, alpha, color, fill, linetype, size
- c + geom_density(kernel = "gaussian")** - x, y, alpha, color, fill, group, linetype, size, weight
- c + geom_dotplot()** - x, y, alpha, color, fill
- c + geom_freqpoly()** - x, y, alpha, color, group, linetype, size
- c + geom_histogram(binwidth = 5)** - x, y, alpha, color, fill, linetype, size, weight
- c2 + geom_qq(aes(sample = hwy))** - x, y, alpha, color, fill, linetype, size, weight

discrete

- d <- ggplot(mpg, aes(f1))
- d + geom_bar()** - x, alpha, color, fill, linetype, size, weight

TWO VARIABLES

continuous x , continuous y

- e <- ggplot(mpg, aes(cty, hwy))
- e + geom_label(aes(label = cty), nudge_x = 1, nudge_y = 1, check_overlap = TRUE)** - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

- e + geom_jitter(height = 2, width = 2)** - x, y, alpha, color, fill, shape, size

- e + geom_point()** - x, y, alpha, color, fill, shape, size, stroke

- e + geom_quantile()** - x, y, alpha, color, group, linetype, size, weight

- e + geom_rug(sides = "bl")** - x, y, alpha, color, linetype, size

- e + geom_smooth(method = lm)** - x, y, alpha, color, fill, group, linetype, size, weight

- e + geom_text(aes(label = cty), nudge_x = 1, nudge_y = 1, check_overlap = TRUE)** - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

discrete x , continuous y

- f <- ggplot(mpg, aes(class, hwy))

- f + geom_col()** - x, y, alpha, color, fill, group, linetype, size

- f + geom_boxplot()** - x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight

- f + geom_dotplot(binaxis = "y", stackdir = "center")** - x, y, alpha, color, fill, group

- f + geom_violin(scale = "area")** - x, y, alpha, color, fill, group, linetype, size, weight

discrete x , discrete y

- g <- ggplot(diamonds, aes(cut, color))

- g + geom_count()** - x, y, alpha, color, fill, shape, size, stroke

THREE VARIABLES

- seals\$z <- with(seals, sqrt(delta_long^2 + delta_lat^2)); l <- ggplot(seals, aes(long, lat))

- l + geom_contour(aes(z = z))** - x, y, z, alpha, colour, group, linetype, size, weight

continuous bivariate distribution

- h <- ggplot(diamonds, aes(carat, price))
- h + geom_bin2d(binwidth = c(0.25, 500))** - x, y, alpha, color, fill, linetype, size, weight

- h + geom_density2d()** - x, y, alpha, colour, group, linetype, size

- h + geom_hex()** - x, y, alpha, colour, fill, size

continuous function

- i <- ggplot(economics, aes(date, unemploy))

- i + geom_area()** - x, y, alpha, color, fill, linetype, size

- i + geom_line()** - x, y, alpha, color, group, linetype, size

- i + geom_step(direction = "hv")** - x, y, alpha, color, group, linetype, size

visualizing error

- df <- data.frame(grp = c("A", "B"), fit = 4.5, se = 1.2)
- j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))

- j + geom_crossbar(fatten = 2)** - x, y, ymax, ymin, alpha, color, fill, group, linetype, size

- j + geom_errorbar()** - x, ymax, ymin, alpha, color, group, linetype, size, width (also **geom_errorbarh()**)

- j + geom_linerange()** - x, ymin, ymax, alpha, color, group, linetype, size

- j + geom_pointrange()** - x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

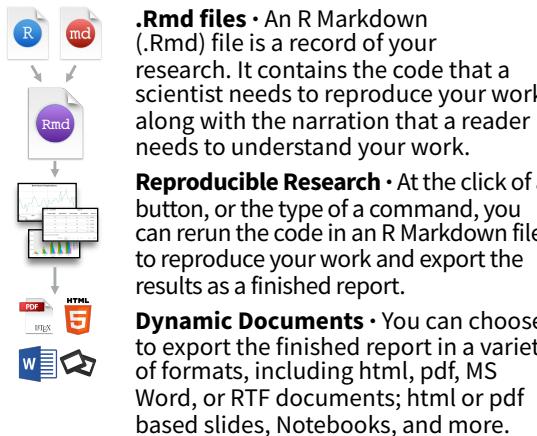
maps

- data <- data.frame(murder = USArrests\$Murder, state = tolower(rownames(USArrests)))
- map <- map_data("state")
- k <- ggplot(data, aes(fill = murder))

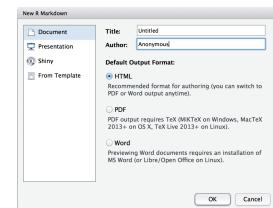
- k + geom_map(aes(map_id = state), map = map)** + **expand_limits(x = map\$long, y = map\$lat)**, map_id, alpha, color, fill, linetype, size

R Markdown :: CHEAT SHEET

What is R Markdown?



Workflow



① Open a new .Rmd file at File ► New File ► R Markdown. Use the wizard that opens to pre-populate the file with a template

② Write document by editing template

③ Knit document to create report; use knit button or render() to knit

④ Preview Output in IDE window

⑤ Publish (optional) to web server

⑥ Examine build log in R Markdown console

⑦ Use output file that is saved along side .Rmd

The screenshot shows the RStudio interface with the following components:

- IDE Area:** Displays the R Markdown file `report.Rmd` with code chunks and their execution status (e.g., `#> [1] '3.2.3'`). Red annotations highlight buttons and menu items: "set preview location" (near the Knit HTML button), "insert code chunk", "run code chunk(s)", "go to code chunk", "publish", "show outline", "modify chunk options", "run all previous chunks", and "run current chunk".
- Output Area:** Shows the rendered HTML output titled "R Markdown" with the content of the R Markdown file. A red annotation points to the "synch publish button to accounts at rpubs.com, shinyapps.io" link.
- Console:** Shows the command `render("report.Rmd", output_file = "report.html")` being run in the R Markdown console.
- File Explorer:** Shows the file structure with `report.Rmd` and `report.html` files.

render

Use `rmarkdown::render()` to render/knit at cmd line. Important args:

input - file to render
output_format

output_options - List of render options (as in YAML)

output_file
output_dir

params - list of params to use

envir - environment to evaluate code chunks in

encoding - of input file

Embed code with knitr syntax

INLINE CODE

Insert with ``r <code>``. Results appear as text without code.

Built with `r getRVersion()` → Built with 3.2.3

CODE CHUNKS

One or more lines surrounded with ````{r}` and `````. Place chunk options within curly braces, after `r`. Insert with `getRVersion()`

```
```{r echo=TRUE}
getRVersion()
```
```

GLOBAL OPTIONS

Set with `knitr::opts_chunk$set()`, e.g.

```
```{r include=FALSE}
knitr::opts_chunk$set(echo = TRUE)
```
```

IMPORTANT CHUNK OPTIONS

cache - cache results for future knits (default = FALSE)

cache.path - directory to save cached results in (default = "cache/")

child - file(s) to knit and then include (default = NULL)

collapse - collapse all output into single block (default = FALSE)

comment - prefix for each line of results (default = "#")

dependson - chunk dependencies for caching (default = NULL)

echo - Display code in output document (default = TRUE)

engine - code language used in chunk (default = 'R')

error - Display error messages in doc (TRUE) or stop render when errors occur (FALSE) (default = FALSE)

eval - Run code in chunk (default = TRUE)

Options not listed above: `R.options`, `aniopts`, `autodep`, `background`, `cache.comments`, `cache.lazy`, `cache.rebuild`, `cache.vars`, `dev`, `dev.args`, `dpi`, `engine.opts`, `engine.path`, `fig.asp`, `fig.env`, `fig.ext`, `fig.keep`, `fig.lp`, `fig.path`, `fig.pos`, `fig.process`, `fig.retina`, `fig.scap`, `fig.show`, `fig.showtext`, `fig.subcap`, `interval`, `out.extra`, `out.height`, `out.width`, `prompt`, `purl`, `ref.label`, `render`, `size`, `split`, `tidy.opts`

fig.align - 'left', 'right', or 'center' (default = 'default')

fig.cap - figure caption as character string (default = NULL)

fig.height, **fig.width** - Dimensions of plots in inches

highlight - highlight source code (default = TRUE)

include - Include chunk in doc after running (default = TRUE)

message - display code messages in document (default = TRUE)

results (default = 'markup')

'asis' - passthrough results

'hide' - do not display results

'hold' - put all results below all code

tidy - tidy code for display (default = FALSE)

warning - display code warnings in document (default = TRUE)

.rmd Structure

rmarkdown

YAML Header

Optional section of render (e.g. pandoc) options written as key:value pairs (YAML).

At start of file

Between lines of ---

Text

Narration formatted with markdown, mixed with:

Code Chunks

Chunks of embedded code. Each chunk:

Begins with ````{r}`

ends with `````

R Markdown will run the code and append the results to the doc.

It will use the location of the .Rmd file as the **working directory**

Parameters

Parameterize your documents to reuse with new inputs (e.g., data, values, etc.)

```
---  
params:  
n: 100  
d: ! Sys.Date()  
---
```

1. **Add parameters** • Create and set parameters in the header as sub-values of params

Today's date is `r params$d`

2. **Call parameters** • Call parameter values in code as `params$<name>`

3. **Set parameters** • Set values with Knit with parameters or the params argument of render():
`render("doc.Rmd", params = list(n = 1, d = as.Date("2015-01-01")))`

Interactive Documents

Turn your report into an interactive Shiny document in 4 steps

1. Add runtime: shiny to the YAML header.
2. Call Shiny input functions to embed input objects.
3. Call Shiny render functions to embed reactive output.
4. Render w `rmarkdown::run` or click Run Document in RStudio IDE

```
---  
output: html_document  
runtime: shiny  
---  
```{r, echo = FALSE}  
numericInput("n",
"How many cars?", 5)
renderTable({
 head(cars, input$n)
})..
```

How many cars?	
5	
1	4.00
2	4.00
3	7.00
4	7.00
5	8.00
	2.00
	10.00
	4.00
	22.00
	16.00

Embed a complete app into your document with `shiny::shinyAppDir()`

**Publish on RStudio Connect**, to share R Markdown documents securely, schedule automatic updates, and interact with parameters in real time.  
[www.rstudio.com/products/connect/](http://www.rstudio.com/products/connect/)





# Pandoc's Markdown

Write with syntax on the left to create effect on right (after render)

```
Plain text
End a line with two spaces
to start a new paragraph.
italics and **bold**
`verbatim` code
sub/superscript22
~~strikethrough~~
escaped: `*` \\
endash: --, emdash: ---
equation: $A = \pi * r^2$
```

```
equation block:
```

```
$$E = mc^2$$
```

```
> block quote
```

```
Header1 {#anchor}
```

```
Header 2 {#css_id}
```

```
Header 3 {.css_class}
```

```
Header 4
```

```
Header 5
```

```
Header 6
```

```
<!--Text comment-->
```

```
\textbf{Text ignored in HTML}
HTML ignored in pdfs
```

```
<http://www.rstudio.com>
[link](www.rstudio.com)
Jump to [Header 1]{#anchor}
image:
```

```
![Caption](smallorb.png)
```

```
* unordered list
+ sub-item 1
+ sub-item 2
- sub-sub-item 1
```

```
* item 2
```

```
Continued (indent 4 spaces)
```

```
1. ordered list
```

```
2. item 2
i) sub-item 1
A. sub-sub-item 1
```

```
(@) A list whose numbering
```

```
continues after
```

```
2. an interruption
```

```
Term 1
```

```
Definition 1
```

Right	Left	Default	Center
12	12	12	12
123	123	123	123
1	1	1	1

```
- slide bullet 1
- slide bullet 2
```

```
(>- to have bullets appear on click)
```

```
horizontal rule/slide break:
```

```

```

```
A footnote [^1]
```

```
[^1]: Here is the footnote.
```

# Set render options with YAML

When you render, R Markdown

1. runs the R code, embeds results and text into .md file with knitr
2. then converts the .md file into the finished format with pandoc



Set a document's default output format in the YAML header:

```

output: html_document

Body
```

**output value**

**creates**

**html\_document**

html

**pdf\_document**

pdf (requires Tex)

**word\_document**

Microsoft Word (.docx)

**odt\_document**

OpenDocument Text

**rtf\_document**

Rich Text Format

**md\_document**

Markdown

**github\_document**

Github compatible markdown

**ioslides\_presentation**

ioslides HTML slides

**slidy\_presentation**

slidy HTML slides

**beamer\_presentation**

Beamer pdf slides (requires Tex)

Customize output with sub-options (listed to the right):

```

output: html_document:
 code_folding: hide
 toc_float: TRUE

Body
```

**html tabs**

Use tablet css class to place sub-headers into tabs

```
Tabset {.tabset .tabset-fade .tabset-pills}
Tab 1
text 1
Tab 2
text 2
End tabset
```



## Create a Reusable Template

1. **Create a new package** with a `inst/rmarkdown/templates` directory

2. In the directory, **Place a folder** that contains:

**template.yaml** (see below)

**skeleton.Rmd** (contents of the template)

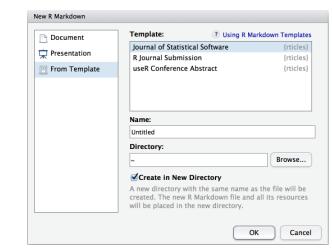
any supporting files

3. **Install the package**

4. **Access template** in wizard at File ▶ New File ▶ R Markdown template.yaml

A footnote <sup>1</sup>

1. Here is the footnote.<sup>1</sup>



sub-option	description	html	pdf	word	odt	rtf	md	gitbook	ioslides	slidy	beamer
<b>citation_package</b>	The LaTeX package to process citations, natbib, biblatex or none	X									
<b>code_folding</b>	Let readers to toggle the display of R code, "none", "hide", or "show"		X								
<b>colortheme</b>	Beamer color theme to use										X
<b>css</b>	CSS file to use to style document		X						X	X	
<b>dev</b>	Graphics device to use for figure output (e.g. "png")	X	X		X	X	X	X	X	X	
<b>duration</b>	Add a countdown timer (in minutes) to footer of slides										X
<b>fig_caption</b>	Should figures be rendered with captions?	X	X	X	X				X	X	X
<b>fig_height, fig_width</b>	Default figure height and width (in inches) for document	X	X	X	X	X	X	X	X	X	X
<b>highlight</b>	Syntax highlighting: "tango", "pygments", "kate", "zenburn", "textmate"	X	X	X					X	X	
<b>includes</b>	File of content to place in document (in_header, before_body, after_body)	X	X	X	X	X	X	X	X	X	
<b>incremental</b>	Should bullets appear one at a time (on presenter mouse clicks)?								X	X	X
<b>keep_md</b>	Save a copy of .md file that contains knitr output	X	X	X	X				X	X	
<b>keep_tex</b>	Save a copy of .tex file that contains knitr output										X
<b>latex_engine</b>	Engine to render latex, "pdflatex", "xelatex", or "lualatex"										X
<b>lib_dir</b>	Directory of dependency files to use (Bootstrap, MathJax, etc.)		X						X	X	
<b>mathjax</b>	Set to local or a URL to use a local/URL version of MathJax to render equations	X							X	X	
<b>md_extensions</b>	Markdown extensions to add to default definition or R Markdown	X	X	X	X	X	X	X	X	X	
<b>number_sections</b>	Add section numbering to headers	X	X								
<b>pandoc_args</b>	Additional arguments to pass to Pandoc	X	X	X	X	X	X	X	X	X	
<b>preserve_yaml</b>	Preserve YAML front matter in final document?										X
<b>reference_docx</b>	docx file whose styles should be copied when producing docx output										X
<b>self_contained</b>	Embed dependencies into the doc										X
<b>slide_level</b>	The lowest heading level that defines individual slides										X
<b>smaller</b>	Use the smaller font size in the presentation?										X
<b>smart</b>	Convert straight quotes to curly, dashes to em-dashes, ... to ellipses, etc.	X								X	X
<b>template</b>	Pandoc template to use when rendering file quarterly_report.html	X	X	X						X	X
<b>theme</b>	Bootswatch or Beamer theme to use for page	X									X
<b>toc</b>	Add a table of contents at start of document	X	X	X	X	X	X	X	X	X	
<b>toc_depth</b>	The lowest level of headings to add to table of contents	X	X	X	X	X	X	X	X	X	
<b>toc_float</b>	Float the table of contents to the left of the main content	X									

## Table Suggestions

Several functions format R data into tables

Table with kable	
eruptions	waiting
3.600	79
1.800	54
3.333	74
2.283	62

eruptionswaiting	
1	3.60
2	1.80
3	3.33
4	2.28

Table with xtable	
eruptions	waiting
3.600	79
1.800	54
3.333	74
2.283	62

```
data <- faithful[1:4,]
```
{r results = 'asis'}
knitr::kable(data, caption = "Table with kable")
```
```
{r results = 'asis'}
print(xtable::xtable(data, caption = "Table with xtable"),
      type = "html", html.table.attributes = "border=0")
```
```
{r results = 'asis'}
stargazer::stargazer(data, type = "html", title = "Table with stargazer")
```

```

&lt;p

## Resources for Future Learning

After this workshop you will have a foundation for building future knowledge and skills in R and data analysis. However, we have barely scratched the surface in terms of what R is, what R can do, and more generally how to code or do data analysis. Below are our favorite resources for learning R.

The most comprehensive inventory of R related resources and educational material can be found on [RStudio's website](#). Check out the huge inventory of guidelines, workshop material, videos, and other useful content.

---

### RStudio Cheatsheets

Very well-designed (if somewhat dense) two-page overview cards for specific R packages or topics. We recommend:

[Data Import Cheatsheet](#): <https://github.com/rstudio/cheatsheets/raw/master/data-import.pdf>

[Data Visualization Cheatsheet](#): <https://github.com/rstudio/cheatsheets/raw/master/data-visualization-2.1.pdf>

[Data Transformation Cheatsheet](#):

<https://github.com/rstudio/cheatsheets/raw/master/data-transformation.pdf>

[R Markdown Cheatsheet](#): <https://github.com/rstudio/cheatsheets/raw/master/rmarkdown-2.0.pdf>

---

### Massive Open Online Courses (MOOCs)/online courses

[The Johns Hopkins University Data Science Specialization](#)

- An excellent series of lessons teaching the A to Z of data science from some of the earliest and best team of R educators.

[Harvard Data Science Specialization](#)

- A fantastic data science course with a leading biostatistician and data scientist.

[stat545](#)

- Data wrangling, exploration, and analysis with R, one of best courses teaching data munging and all things R, initially taught by statistician Jenny Bryan at UBC.
- 

### Online/Free Textbooks

[R for Data Science by Hadley Wickham and Garrett Grolemund](#)

- The definitive text on data science via the tidyverse by the leading R developer Hadley Wickham.

[Introduction to Data Science by Rafael Irizarry](#)

- Raf Irizarry, Harvard Professor and fantastic teacher has published a wonderful introductory Data Science Book.

[An Introduction to Statistical and Data Sciences via R by Chester Ismay and Albert Y. Kim](#)

- A fantastic introduction to R via statistical inference.

[Fundamentals of Data Visualization by Clause Wilke](#)

- Claus Wilke, a professor from UT Austin has written a highly useful ggplot [add-on package](#) and now has a new book on data visualization.
- 

## **Youtube videos**

Anything with Hadley Wickham, including:

- [Hadley Wickham's "dplyr" tutorial at useR 2014](#)
- [Stanford Seminar - Expressing yourself in R](#)

As well as Garrett Grolemund's [Data Wrangling Series](#)

---

## **Websites/Blogs**

[R Bloggers](#) is a blog aggregator which captures a lot of the most prominent R bloggers on the internet

[Rweekly.org](#) is a weekly manual review of the latest and greatest in R internet content.

[Rviews](#) is the RStudio blog devoted to the R Community and the R Language.