

How to Truly "Excel" at Data Analysis and Visualization: An Introduction to the R Programming Language

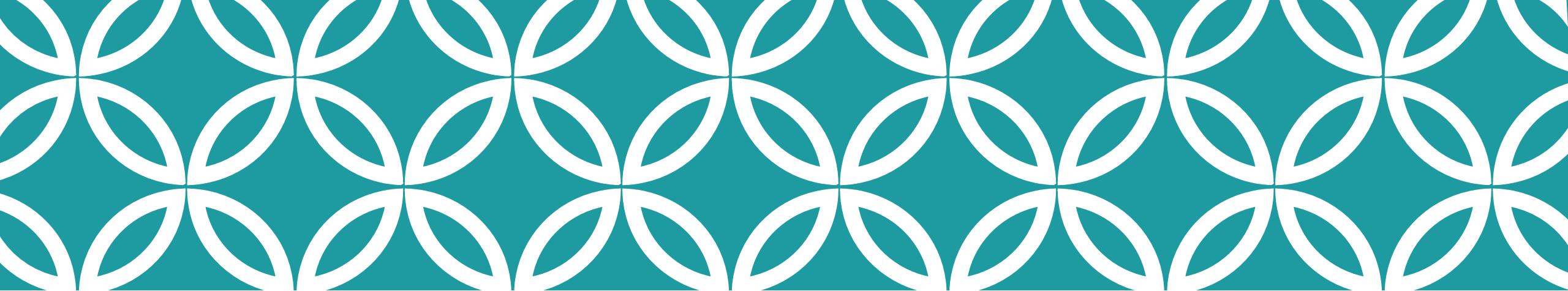
**Daniel Herman
Patrick Mathias
Joseph Rudolf
December 13, 2020**

Course Introduction

Course Goals

1. Gain familiarity with the R programming language and RStudio software
2. Learn some fundamentals of coding using Tidyverse tools on a clinical data set
3. Develop workflows for reproducible data analysis using RMarkdown

December 13 2020	Session	Instructor
8:30 am - 8:50 am	Introduction to the Course	Daniel Herman Patrick Mathias Joseph Rudolf
9:00 am - 9:45 am	Introduction to R and RStudio	Joseph Rudolf
10:00 am - 10:45 am	Reproducible Reporting	Patrick Mathias
11:00 am - 11:00 am	Data Visualization	Joseph Rudolf
12:30 pm - 1:15 pm	Data Slicing	Patrick Mathias
1:30 pm - 2:15 pm	Data Transformation	Daniel Herman
2:30 pm - 3:15 pm	Statistical Analysis	Daniel Herman



Who are we?

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



Patrick Mathias

Assistant Professor, Department of
Laboratory Medicine

University of Washington School of
Medicine

Associate Medical Director, Laboratory
Medicine Informatics



Joseph Rudolf

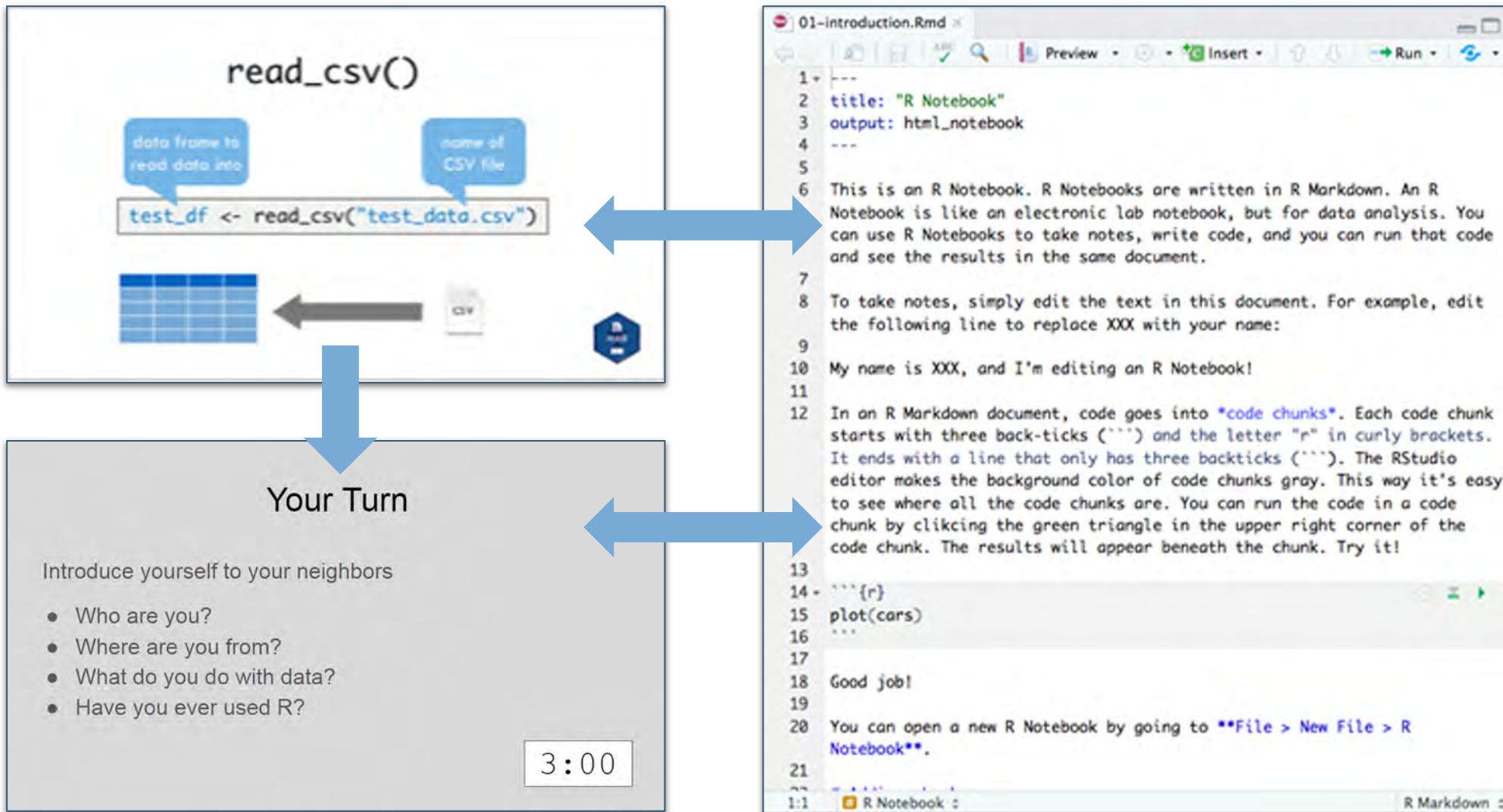
Assistant Professor, Department of Pathology,
University of Utah Medical School

Medical Director, Automated Core Laboratory, ARUP
Laboratories



Workshop Workflow

Sessions

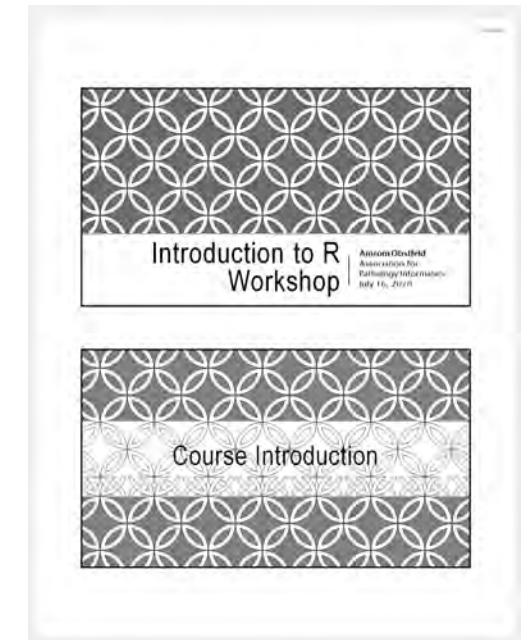


Configuring Your Setup

A screenshot of the RStudio interface. The top menu bar includes File, Edit, Code, View, Project, Workspace, Plots, Tools, and Help. The left sidebar shows a project named "diamondPricing.Rproj" with files "formatPlot.R" and "diamonds.R". The main workspace shows the "diamonds" dataset with 53940 observations and 10 variables. A summary of the "price" variable is displayed. The bottom panel shows a scatter plot titled "Diamond Pricing" with "Carat" on the x-axis and "Price" on the y-axis. The plot is color-coded by "Clarity" levels, with points ranging from small (I1) to large (IF).

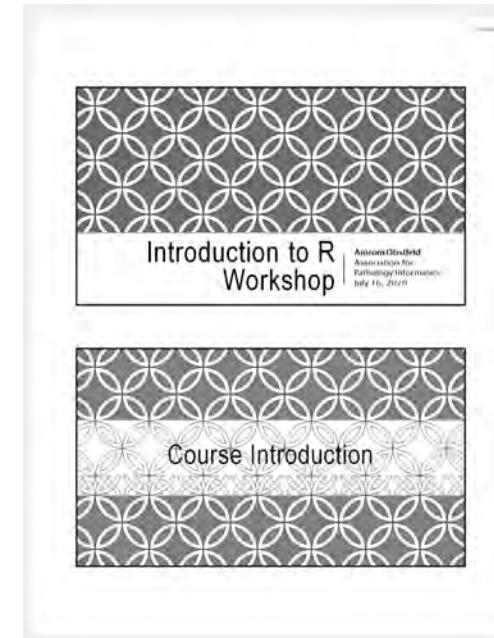
```
library(ggplot2)
source("plots/formatPlot.R")
view(diamonds)
summary(diamonds)
p <- qplot(carat, price,
           data=diamonds, color=clarity,
           xlab="carat", ylab="price",
           main="Diamond Pricing")
format.plot(p, size=24)
```

```
Min. : 0.000 1st Qu.: 4.710  Median : 5.700  Mean  : 5.731 3rd Qu.: 6.540  Max. :10.740
1st Qu.: 0.000 1st Qu.: 4.720 1st Qu.: 2.910 Median : 5.700 Median : 5.710 Median : 3.530 Mean   : 5.731 Mean   : 5.735 Mean   : 3.539 3rd Qu.: 6.540 3rd Qu.: 6.540 3rd Qu.: 4.040 Max.  :10.740 Max.  :58.900 Max.  :31.800
```



Workshop Coursebook

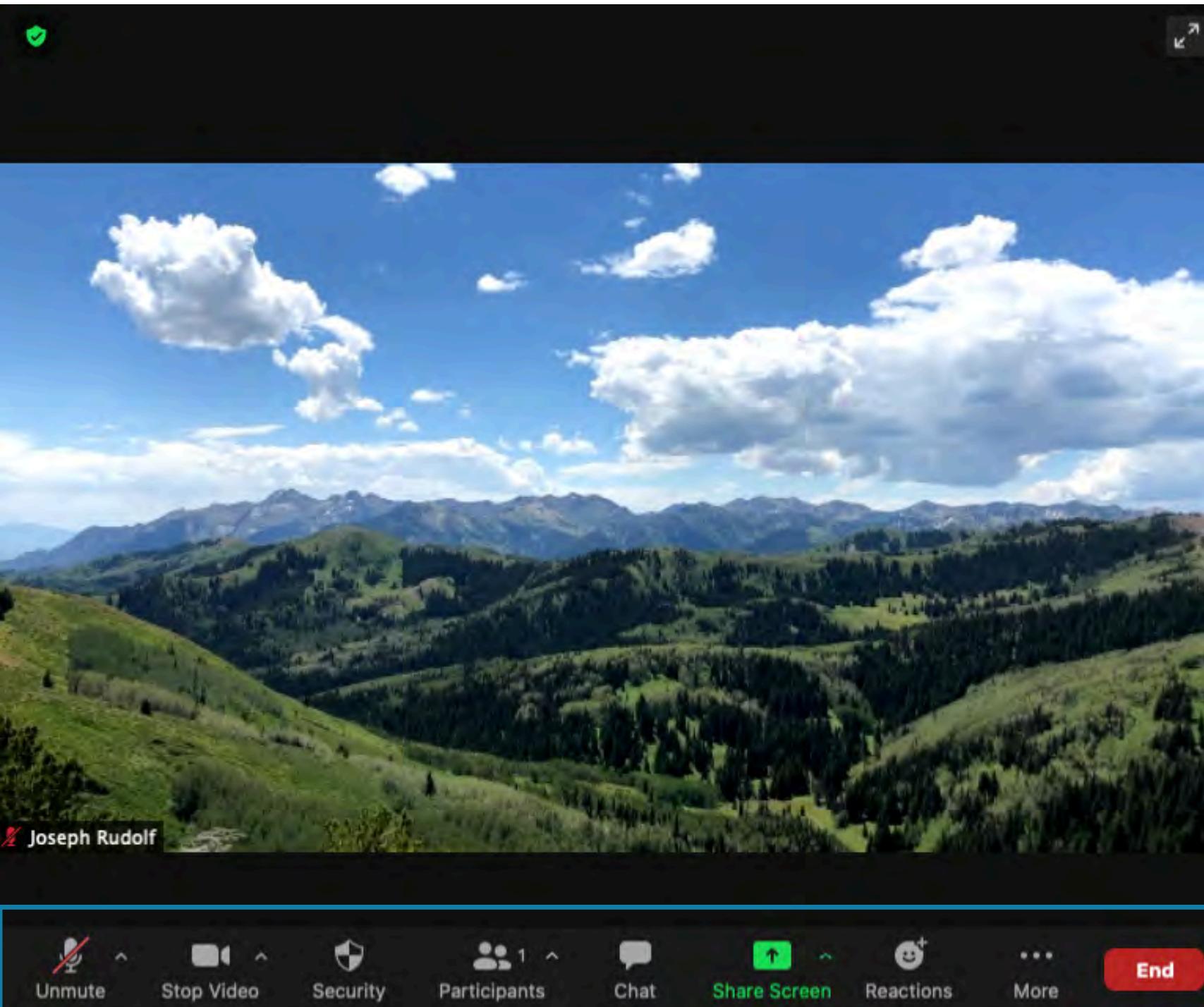
- Print out of all slides
- Appendix
 - Cheat Sheets
 - Useful resources



Using Zoom



- Participants muted
- Chat window
- Non-verbal feedback



Participants (1)

JR Joseph Rudolf (Host, me) X End

yes no go slower go faster more clear all

Invite Mute All More

Chat

To: Everyone File ...

Type message here...

Getting Help

- During presentation – Enter your question in the chat, an instructor will respond to the group or direct message you as appropriate.
- Between sessions – Feel free to unmute yourself and converse with instructors.



Who are you?

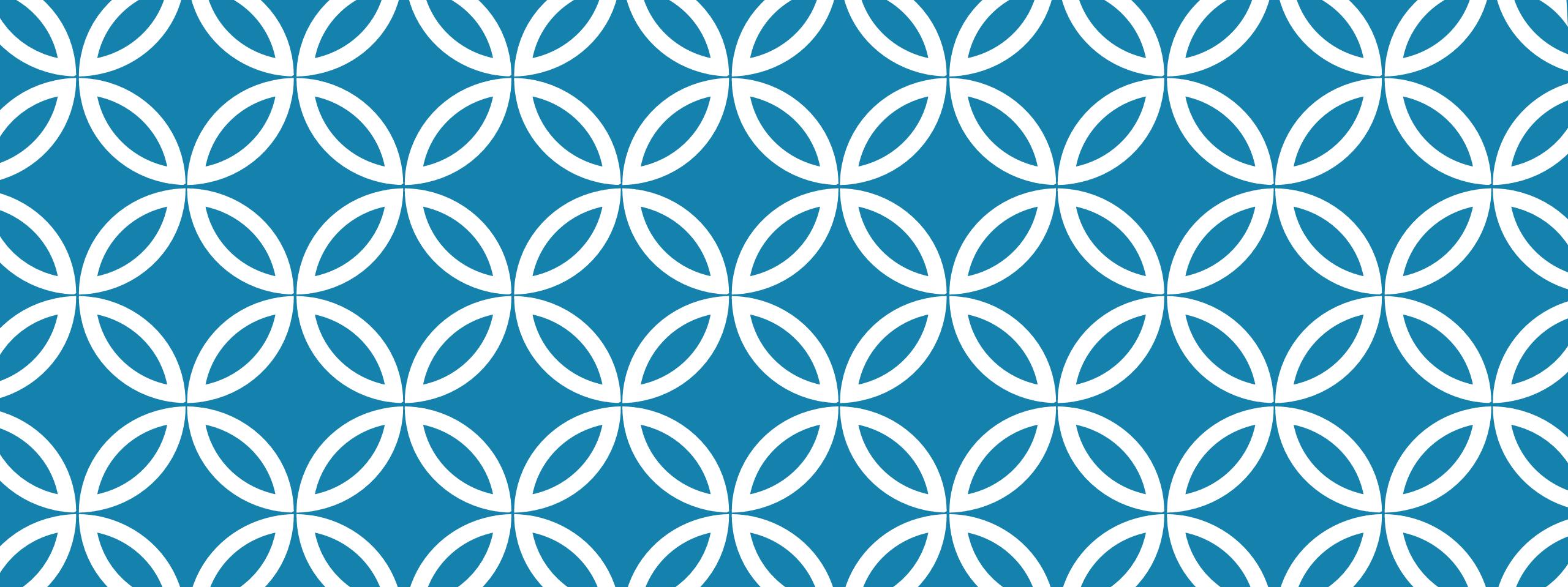
Your Turn

Where are you today?

Type your name and location in the chat.

Pre-Course Survey Results

Time for the Course!!



Introduction to R and R Studio

Session 1
Joseph Rudolf
December 13, 2020

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: <https://tinyurl.com/r-aacc-2020>

Enter your log in credentials

Join Space

Make a copy of the Core Exercises for yourself

Join Space?

Joining a space gives you access to it and to its contents.

Once you join, admins will be able to see your email address.

Would you like to join this space?

Join Space

Cancel



https://rstudio.cloud/spaces/110165/projects

AACC Introduction to R Course Projects Members About Joseph Rudolf

All Projects New Project

List All projects Sort By name

AACC-2020-Introduction-to-R Patrick Mathias Created Dec 9, 2020 4:28 PM

Copy Delete Move

1

2

3

Copy Project

Make your own copy of AACC-2020-Introduction-to-R?

OK

R Studio Cloud Terms Status

in f t q © 2020 RStudio, PBC

This screenshot shows the RStudio Cloud interface for managing projects. On the left, a sidebar lists various sections like 'Spaces', 'Learn', 'Help', and 'Info'. The main area displays a list of projects under 'All Projects'. One specific project, 'AACC-2020-Introduction-to-R', is highlighted with a blue box. In the top right of this project card, there are buttons for 'Copy', 'Delete', and 'Move'. Three numbered arrows with labels point to these actions: '1' points to the 'New Space' button in the sidebar; '2' points to the 'Copy' button on the project card; and '3' points to the 'OK' button in a modal dialog titled 'Copy Project' that has popped up over the project card.

R Studio Cloud  AACC Introduction To R Course / AACC-2020-Introduction-to-R

File Edit Code View Plots Session Build Debug Profile Tools Help

Spaces Your Workspace AACC Introduction to R Courses New Space

Console Terminal Jobs /cloud/project/

R version 4.0.3 (2020-10-10) -- "Bunny-Wunnies Freak Out"
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.

Learn Guide What's New Primers Cheat Sheets Help Current System Status RStudio Community

Info Plans & Pricing Terms and Conditions

Environment History Connections Git Tutorial Import Dataset Global Environment

Environment is empty

Files Plots Packages Help Viewer

New Folder Upload Delete Rename More

Cloud project

	Name	Size	Modified
<input type="checkbox"/>	..		
<input type="checkbox"/>	.gitignore	621 B	Dec 9, 2020, 4:29 PM
<input type="checkbox"/>	Rhistory	0 B	Dec 9, 2020, 4:29 PM
<input type="checkbox"/>	03 - Visualize.Rmd	3 KB	Dec 9, 2020, 4:29 PM
<input type="checkbox"/>	04 - Slice.Rmd	3.8 KB	Dec 10, 2020, 11:19 PM
<input type="checkbox"/>	04 - Transform - b.Rmd	2.3 KB	Dec 10, 2020, 10:22 AM
<input type="checkbox"/>	04 - Transform.Rmd	6.1 KB	Dec 9, 2020, 4:29 PM
<input type="checkbox"/>	05 - Stats.Rmd	4.6 KB	Dec 9, 2020, 4:29 PM
<input type="checkbox"/>	coursepack	205 B	Dec 12, 2020, 12:03 PM
<input type="checkbox"/>	data	6.6 KB	Dec 10, 2020, 11:19 PM
<input type="checkbox"/>	presentations	3.6 KB	Dec 9, 2020, 4:29 PM
<input type="checkbox"/>	project.Rproj		
<input type="checkbox"/>	README.md		
<input type="checkbox"/>	resources.md		
<input type="checkbox"/>	src		

Your Turn

Navigate to: <https://tinyurl.com/r-aacc-2020>

Enter your log in credentials

Join Space

Make a copy of the Core Exercises for yourself

Console Terminal Jobs

```
/cloud/project/ >

R version 3.5.2 (2018-12-20) -- "Eggshell Igloo"
Copyright (C) 2018 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)

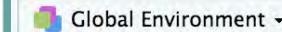
R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

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



List



Environment is empty

ENVIRONMENT

Files Plots Packages Help Viewer

New Folder Upload Delete Rename More

Cloud > project

	Name	Size	Modified
	..		
	.Rhistory	0 B	Apr 20, 2019, 1:39 PM
	project.Rproj	205 B	Apr 20, 2019, 1:39 PM

OUTPUT

CONSOLE

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:

- Take the integer 1974
- Subtract 12
- Multiply by 29

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 Basics of Coding: Functions

- Code that extends our reach beyond the basic operators

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

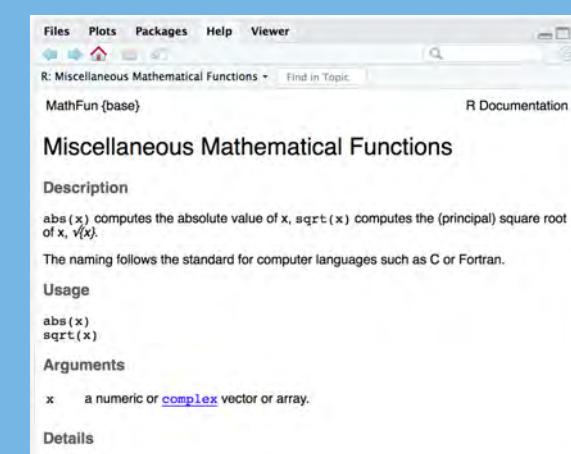
- What if I don't know what a function does?

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

function
(does stuff)

argument
(input)

abs(-77)



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/>)

Putting Functions to Work

- We can use functions to do more than simple math, we can make things!
- We can create a series of integers (a vector) using the `seq()` function

```
>  
> seq(from=5, to=150, by=10)  
[1] 5 15 25 35 45 55 65 75 85 95 105 115 125 135 145
```

The Basics of Coding: Objects

- Objects are the container for your output

object

(stores output)

function

(does stuff)

arguments

(input)

```
sequence_of_10s <- seq(from=5, to=150, by=10 )
```

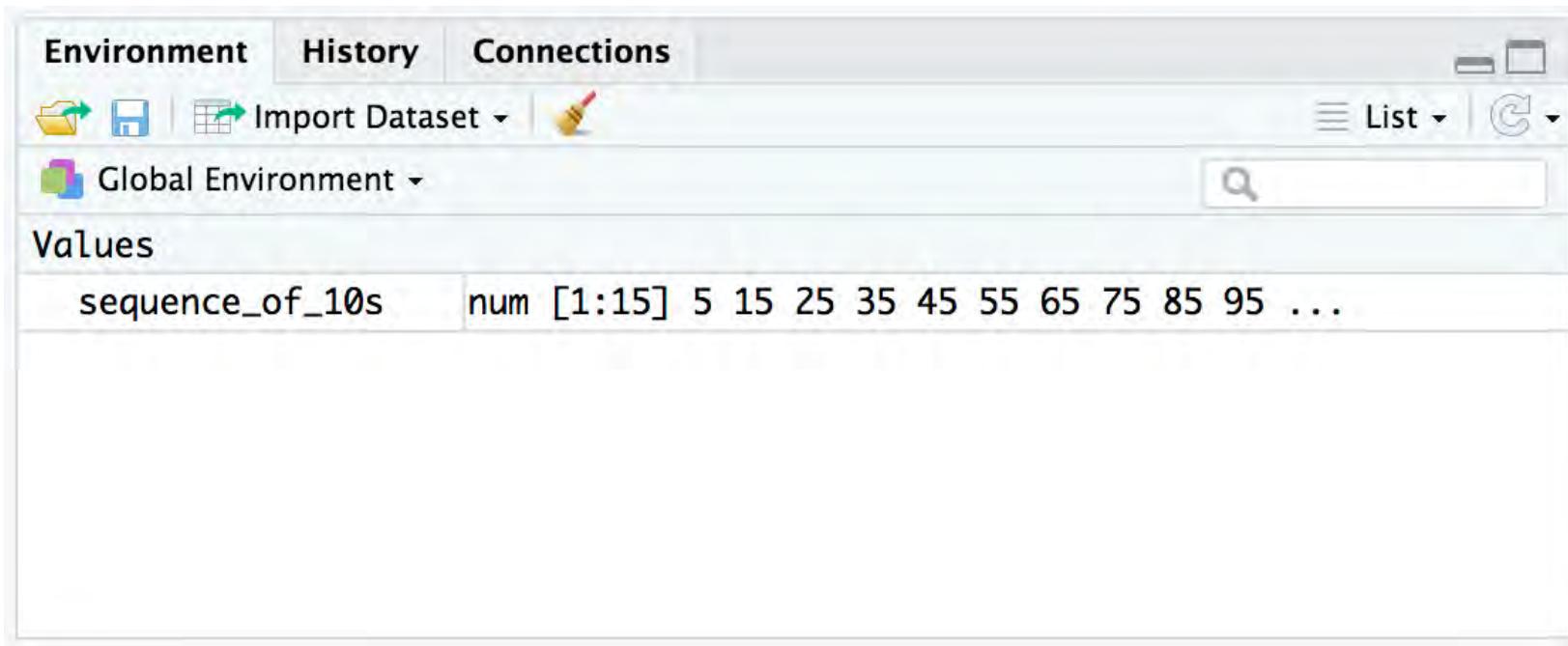
Checking the Contents of an Object

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

```
>  
> sequence_of_10s  
[1] 5 15 25 35 45 55 65 75 85 95 105 115 125 135 145
```

Checking the contents of an object

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



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!

```
>  
> min(sequence_of_10s)  
[1] 5  
>  
> max(sequence_of_10s)  
[1] 145  
>
```

Your Turn 2

Generate a sequence, store it to an object, and **ply** your object

Type the following code to create a sequence from 0 to 500 in increments of 25 called `sequence_of_25s`:

```
sequence_of_25s <- seq(from=0, to=500, by=25)
```

Calculate the median value of this series using the `median()` function

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.



Installing and Loading Packages

- Installing a package

function
(does stuff)

arguments
(input)

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

- Loading into your environment

```
library(tidyverse)
```

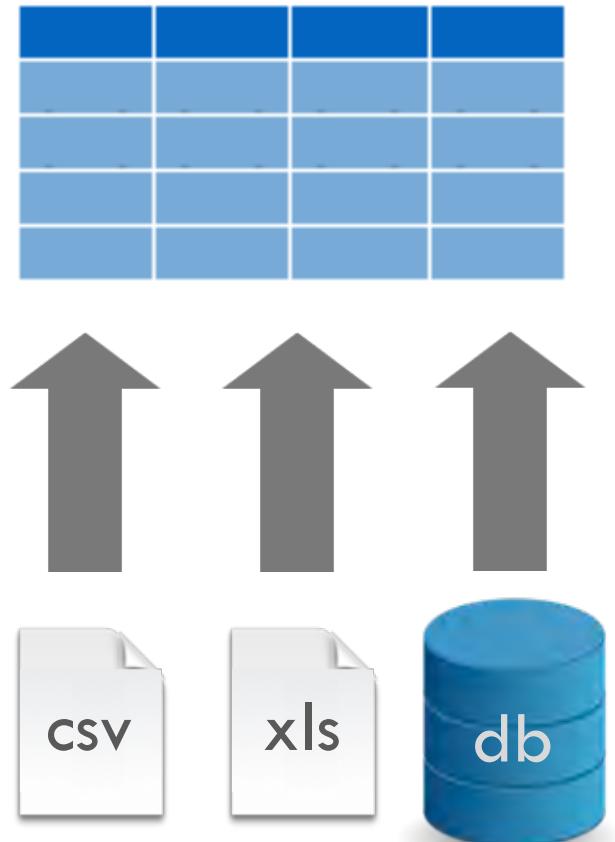


Importing Data and Working with Dataframes (aka Useful Data)



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

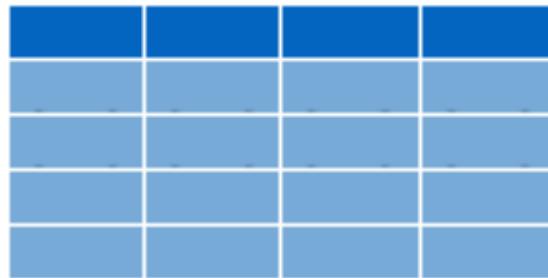
rectangular
structure

The screenshot shows a Microsoft Excel window with the title bar '02-example'. The spreadsheet contains four rows of data. The first row, 'Name,MRN,DOB', serves as a header. The subsequent three rows provide data points: 'Santa Claus,12345,1/1/01', 'Roger Rabbit,67890,12/12/69', and 'Kermit the Frog,24680,2/2/22'. The data is presented in a simple, rectangular, flat file structure.

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

Loading Data to Create a Dataframe

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



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`

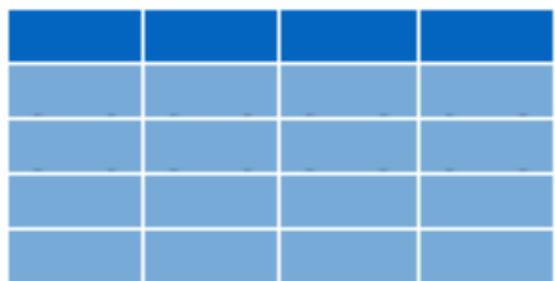
`read_csv()`

data frame to
read data into

name of
CSV file

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

`covid_testing`



`covid_testing.csv`



What's in a name?

- Capitalization matters

covid_testing ≠ Covid_testing ≠ COVID_TESTING

- 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

15,524
Observations
(Rows)

	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	negative
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

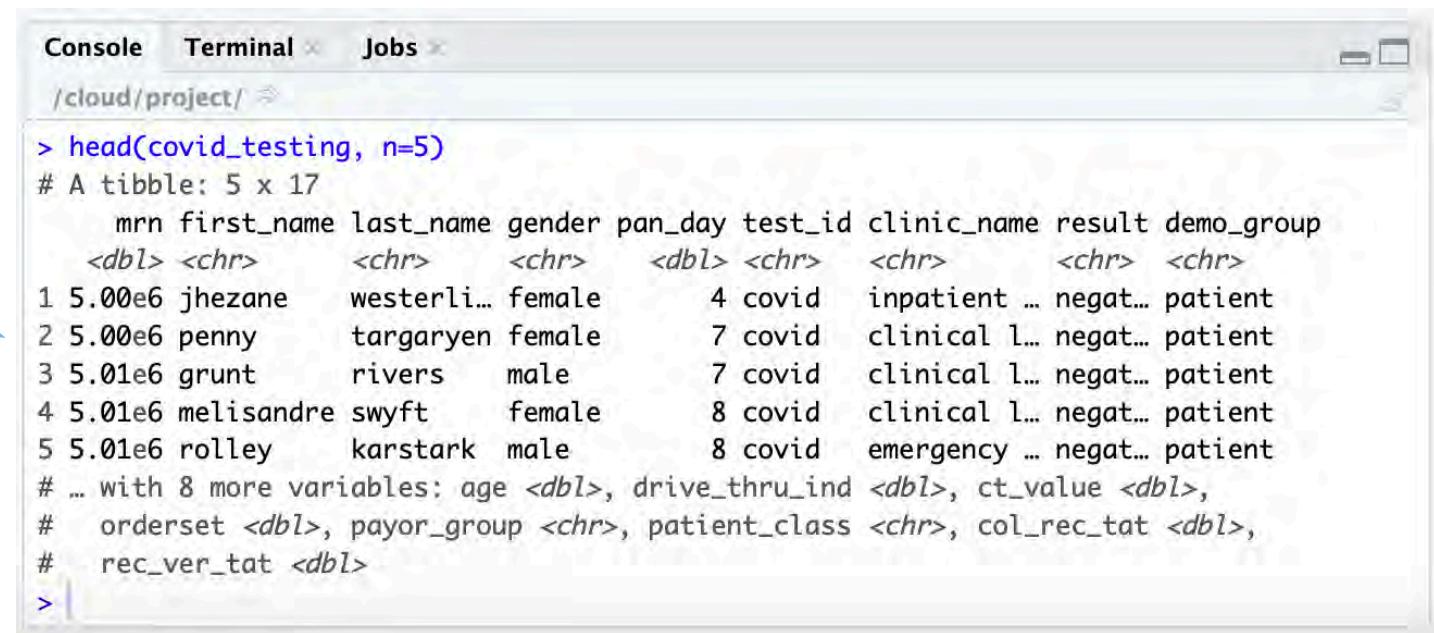
17 Attributes
(Columns)

Working with Dataframes at the Console

- The `head()` function is helpful for displaying a snippet of your dataframe

```
head(object_name, n=number of rows to view)
```

Sample of Data
in Your Object



A screenshot of the RStudio interface showing the Console tab. The command `> head(covid_testing, n=5)` is entered, followed by the resulting output which shows the first 5 rows of a tibble named `covid_testing`. The columns are labeled: mrn, first_name, last_name, gender, pan_day, test_id, clinic_name, result, demo_group, and several unnamed columns represented by '<dbl>' and '<chr>'. The data includes patient information such as names like jhezane, penny, grunt, melisandre, and roolley, along with their medical details like gender, test results, and clinical status.

```
Console Terminal × Jobs ×  
/cloud/project/  
> head(covid_testing, n=5)  
# A tibble: 5 x 17  
   mrn first_name last_name gender pan_day test_id clinic_name result demo_group  
   <dbl> <chr>     <chr>    <chr>    <dbl> <chr>     <chr>    <chr>    <chr>  
1 5.00e6 jhezane  westerli... female      4 covid  inpatient ... negat... patient  
2 5.00e6 penny    targaryen female      7 covid  clinical l... negat... patient  
3 5.01e6 grunt    rivers    male       7 covid  clinical l... negat... patient  
4 5.01e6 melisandre swyft    female      8 covid  clinical l... negat... patient  
5 5.01e6 rolley   karstark   male       8 covid  emergency ... negat... patient  
# ... with 8 more variables: age <dbl>, drive_thru_ind <dbl>, ct_value <dbl>,  
#   orderset <dbl>, payor_group <chr>, patient_class <chr>, col_rec_tat <dbl>,  
#   rec_ver_tat <dbl>  
>
```

Your Turn 4

Understanding Object Contents

Use the `tail()` function to view the last 10 rows of our object `covid_testing`.

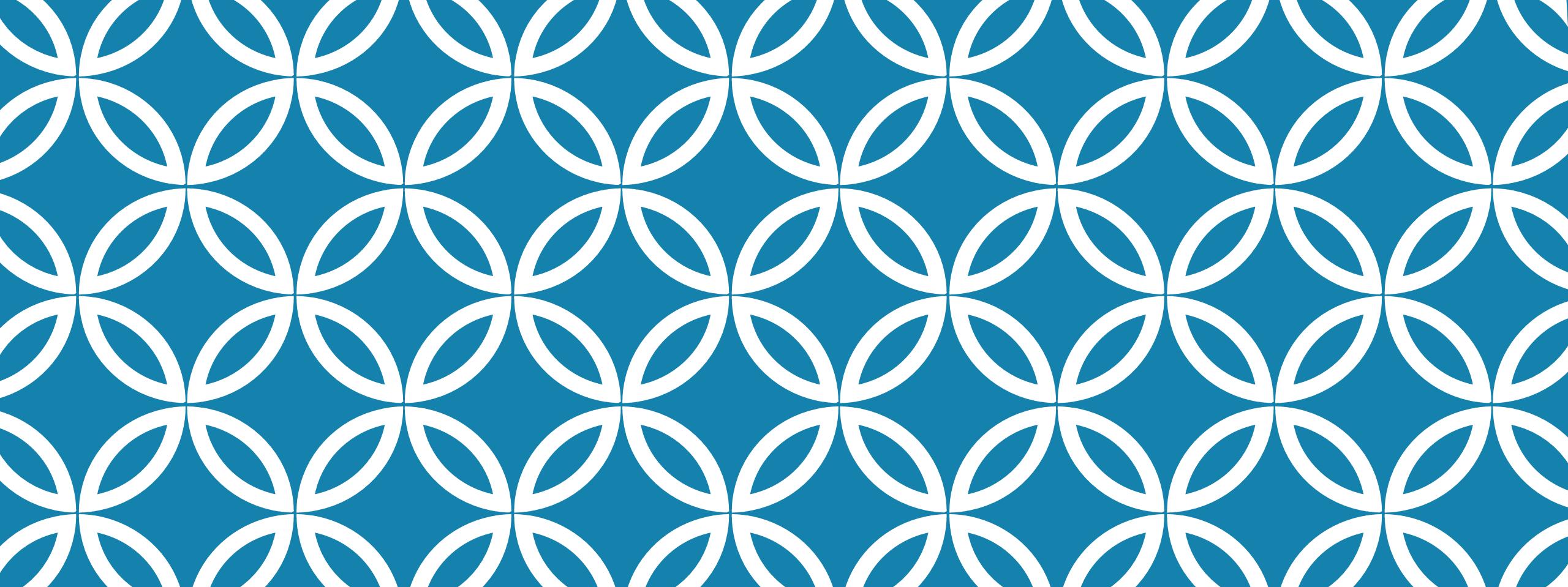
-What is the ratio of female to male patients in this subset of data?

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

Patrick Mathias
December 13, 2020

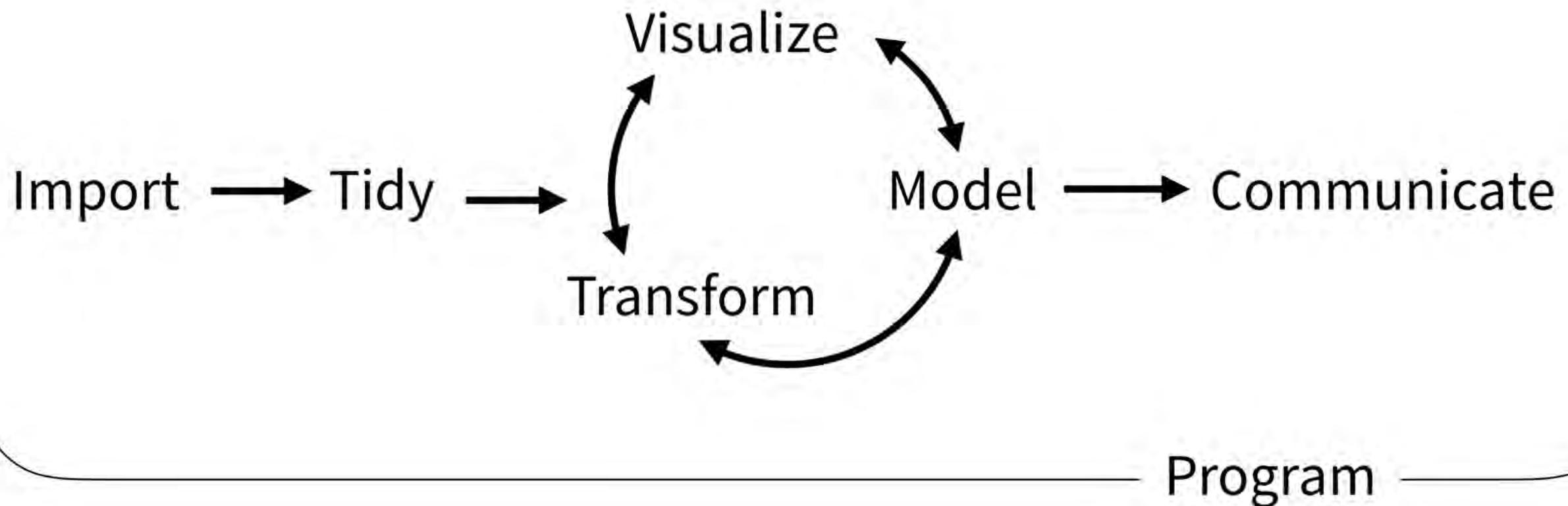
Goals

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

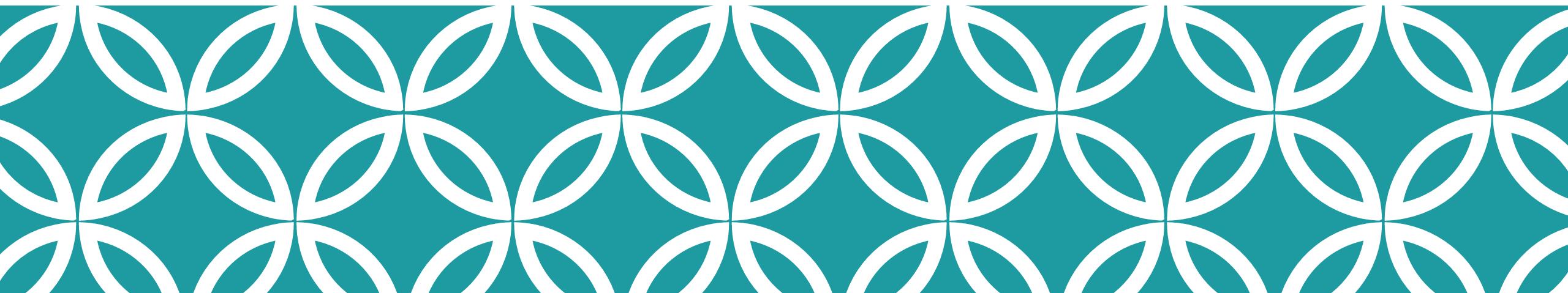
Typical Data Science Pipeline



Your Turn #1

1. List 3 benefits of doing data analysis using spreadsheet software like Microsoft Excel
2. List 3 drawbacks/problems with doing data analysis using spreadsheet software like Microsoft Excel





Why is reproducibility important?

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

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)**

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



Using R Markdown to Support Reproducibility



```
File Edit Code View Plots Session Build Debug Profile Tools Help
+ | Go to file/function Addins
04 - Stats.Rmd
| ABC Preview | Insert | Run |
1 ---  
2 title: "Summarization and statistics in R"  
3 output: html_notebook  
4 editor_options:  
5   chunk_output_type: inline  
6 ---  
7  
8 * ## Setup  
9  
10 * `r`{r setup}  
11   library(tidyverse)  
12   library(readxl)  
13  
14   orders <- read_excel("data/orders_data_set.xlsx")  
15  
16  
17 * ## Summarize  
18  
19   `r`{r}  
20   orders %>%  
21     select(order_id, patient_id) %>%  
22     head(4) %>%  
23     summarize(order_count = n(),  
24               pt_count    = n_distinct(patient_id))  
25  
26  
27 * ## Your Turn 1  
28  
29   Add onto the code in the above chunk to calculate:  
30  
31   1) Mean count of orders per patient
```

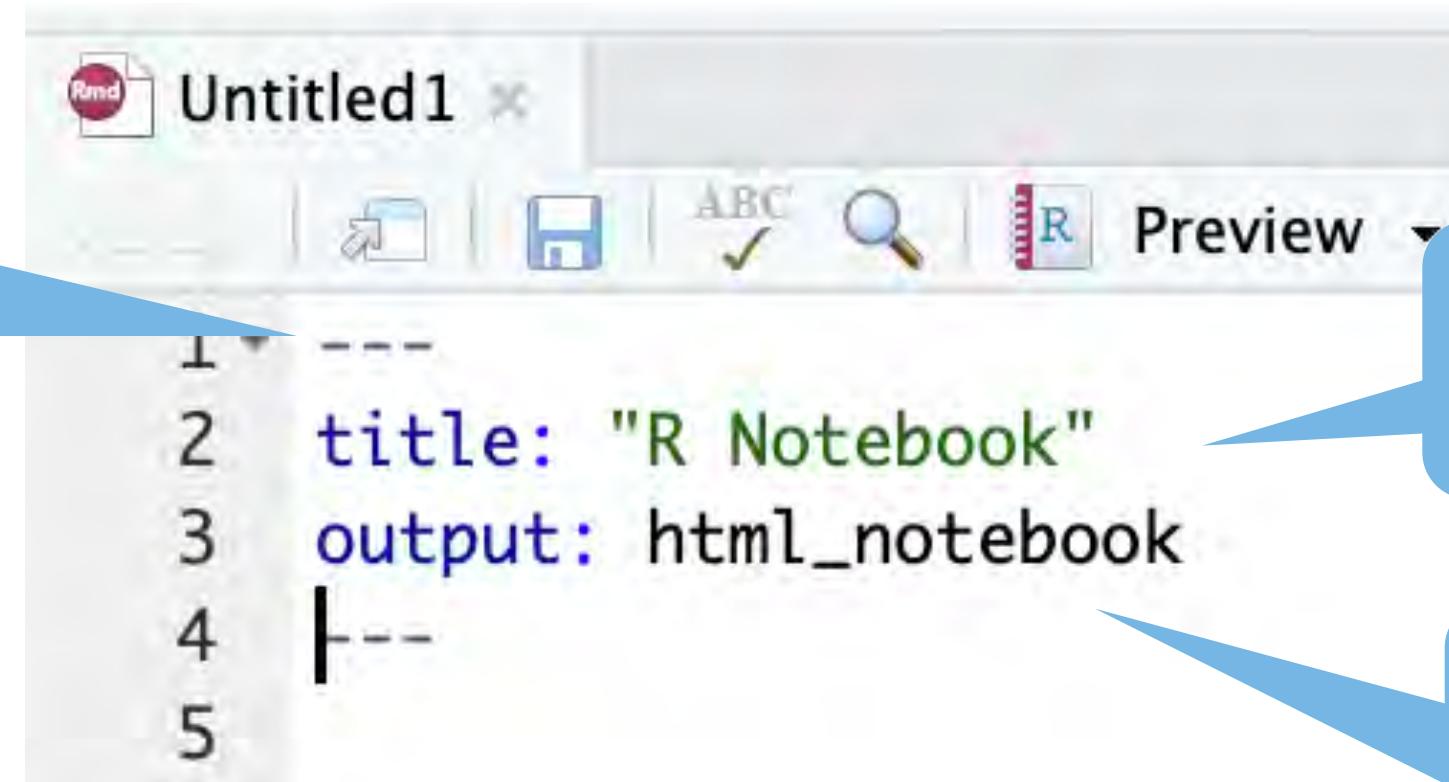
Header

Code chunk

Text

Header

Starts and ends
with 3 dashes



A screenshot of an R Notebook interface titled "Untitled1". The notebook contains five code cells:

```
1 ---  
2 title: "R Notebook"  
3 output: html_notebook  
4 ---  
5
```

Title in quotes

Output format

Text

5
6 This is an [R Markdown](<http://rmarkdown.rstudio.com>) Notebook. When you execute
code within the notebook, the results appear beneath the code.

7
8 Try executing this chunk by clicking the *Run* button within the chunk or by
placing your cursor inside it and pressing *Cmd+Shift+Enter*.

9

1 asterisk for *italics* (*italics*)

2 asterisks for **bold** (**bold**)

Hyphens (-bullet 1) for bullet points

Include hyperlinks
with [text](link)
format

1 asterisk for
italics (*italics*)

R Markdown

Cheat Sheet
learn more at rmarkdown.rstudio.com

markdown 0.2.50 Updated: 8/14



2. Open File

Start by saving a text file with the extension .Rmd, or open an RStudio Rmd template

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
##### Header 5
##### Header 6

endash: --
emdash: ---
ellipsis: ...
inline equation: $A = \pi * r^2$
image: 

horizontal rule (or slide break):
```

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
Header 4
Header 5
Header 6

endash: --
emdash: ---
ellipsis: ...
inline equation: $A = \pi * r^2$
image:

horizontal rule (or slide break):

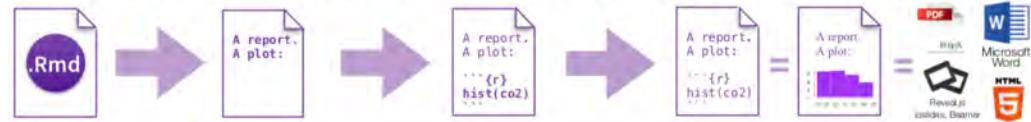
output: beamer_presentation..... beamer slideshow (pdf)
output: ioslides_presentation..... ioslides slideshow (html)

RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com

1. Workflow

R Markdown is a format for writing reproducible, dynamic reports with R. Use it to embed R code and results into slideshows, pdfs, html documents, Word files and more. To make a report:

- i. **Open** - Open a file that uses the .Rmd extension.
- ii. **Write** - Write content with the easy to use R Markdown syntax
- iii. **Embed** - Embed R code that creates output to include in the report
- iv. **Render** - Replace R code with its output and transform the report into a slideshow, pdf, html or ms Word file.



3. Markdown

Next, write your report in plain text. Use markdown syntax to describe how to format text in the final report.

syntax

Plain text
End a line with two spaces to start a new paragraph.
italics and *italics*
bold and **bold**

> block quote

* unordered list

* item 2

- + sub-item 1
- + sub-item 2

1. ordered list

2. item 2

- + sub-item 1
- + sub-item 2

Table Header | Second Header

Table Header	Second Header
Table Cell	Cell 2
Cell 3	Cell 4

1. ordered list

2. item 2

- + sub-item 1
- + sub-item 2

Table Header | Second Header

Table Header	Second Header
Table Cell	Cell 2
Cell 3	Cell 4

becomes

Plain text
End a line with two spaces to start a new paragraph.
italics and *italics*
bold and **bold**

horizontal rule (or slide break):

block quote

- unordered list
- item 2
 - sub-item 1
 - sub-item 2

1. ordered list
2. item 2
 - sub-item 1
 - sub-item 2

Table Header Second Header

Table Header	Second Header
Table Cell	Cell 2
Cell 3	Cell 4

1. ordered list

2. item 2

- sub-item 1
- sub-item 2

Table Header Second Header

Table Header	Second Header
Table Cell	Cell 2
Cell 3	Cell 4

Your Turn #2

1. Open up an R Notebook per guidance on next few slides.
2. Change the title of your R Notebook to “My First R Notebook” by modifying the header.
3. Add your name as the author by adding another line to the header:
author: “Your Name”
3. Add a second level heading (##) at the end of the notebook called “My Calculation”



Step 1

The screenshot shows the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, and Help. Below the menu is a toolbar with icons for file operations like New, Open, Save, and Print. A search bar says "Go to file/function" and an "Addins" dropdown is shown. The main workspace shows a file named "04-26.R" with the title "Planting of a Tree". The code in the file is:

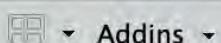
```
## -- "Planting of a Tree"
## R Foundation for Statistical Computing
## nux-gnu (64-bit)
```

Below the code, two annotations are present: "Step 1" with a blue arrow pointing to the "New" icon in the toolbar, and "Step 2" with a blue arrow pointing to the "R Notebook" option in the "File" menu.

The right side of the interface features several panes: an "Environment" pane listing Global Environment, Data, and test_catalog; a "Files" pane showing a directory tree with files like .Rhists, .Rprof, 02 - R, 03 - ti, 04 - S, 05 - V, 06 - E, data, and project; and a "Plots" pane which is currently empty.



Go to file/function



Addins

Console Terminal Jobs

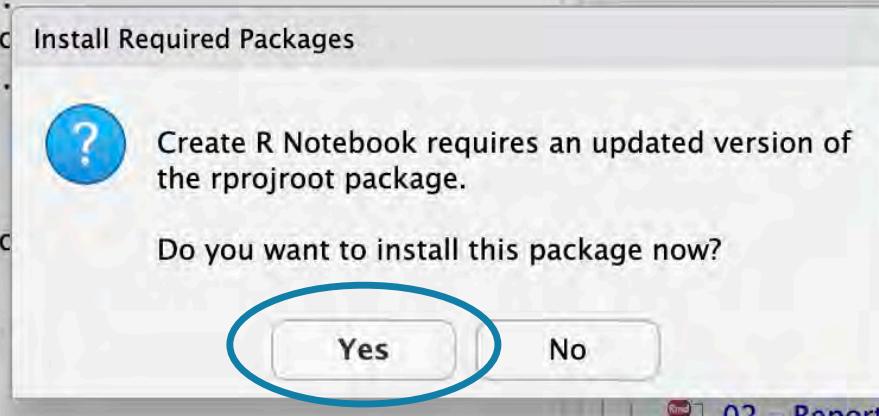
/cloud/project/

R version 3.6.0 (2019-04-26) -- "Planting of a Tree"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)

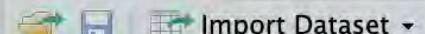
R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

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,
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.



Environment History Connections



Global Environment

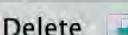
Data

test_catalog

559 obs. of 2 variables

Pages Help Viewer

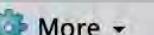
Upload



Delete



Rename



		Size	Modified
	02 – Report.Rmd	0 B	Jul 4,
	03 – transform.Rmd	69 B	Jul 4,
	04 – Stats.Rmd	2.4 KB	Jul 4,
	05 – Visualize.Rmd	3.6 KB	Jul 4,
	06 – Exploratory Data Analysis.Rmd	2.7 KB	Jul 4,
	data	9.9 KB	Jul 4,
	project.Rproj	205 B	Jul 4,

You may see this screen – click yes

+ | Go to file/function | Addins

Untitled1

1 ---
2 title: "R Notebook"
3 output: html_notebook
4 ---
5
6 This is an [R Markdown](http://rmarkdown.rstudio.com) Notebook. When you execute code within the notebook, the results appear beneath the code.
7
8 Try executing this chunk by clicking the *Run* button within the chunk or by placing your cursor inside it and pressing *Cmd+Shift+Enter*.
9
10 ``{r}
11 plot(cars)
12 ``
13
14 Add a new chunk by clicking the *Insert Chunk* button on the toolbar or by pressing *Cmd+Option+I*.
15
16 When you save the notebook, an HTML file containing the code and output will be saved alongside it (click the *Preview* button or press *Cmd+Shift+K* to preview the HTML file).
17
18 The preview shows you a rendered HTML copy of the contents of the editor.

4:1 # R Notebook

Your Turn #2

1. Open up an R Notebook per guidance on next few slides.
2. Change the title of your R Notebook to “My First R Notebook” by modifying the header.
3. Add your name as the author by adding another line to the header:
Author: “Your Name”
3. Add a second level heading (##) at the end of the notebook called “My Calculation”



Code chunks

Open/close with
3 backticks

Language

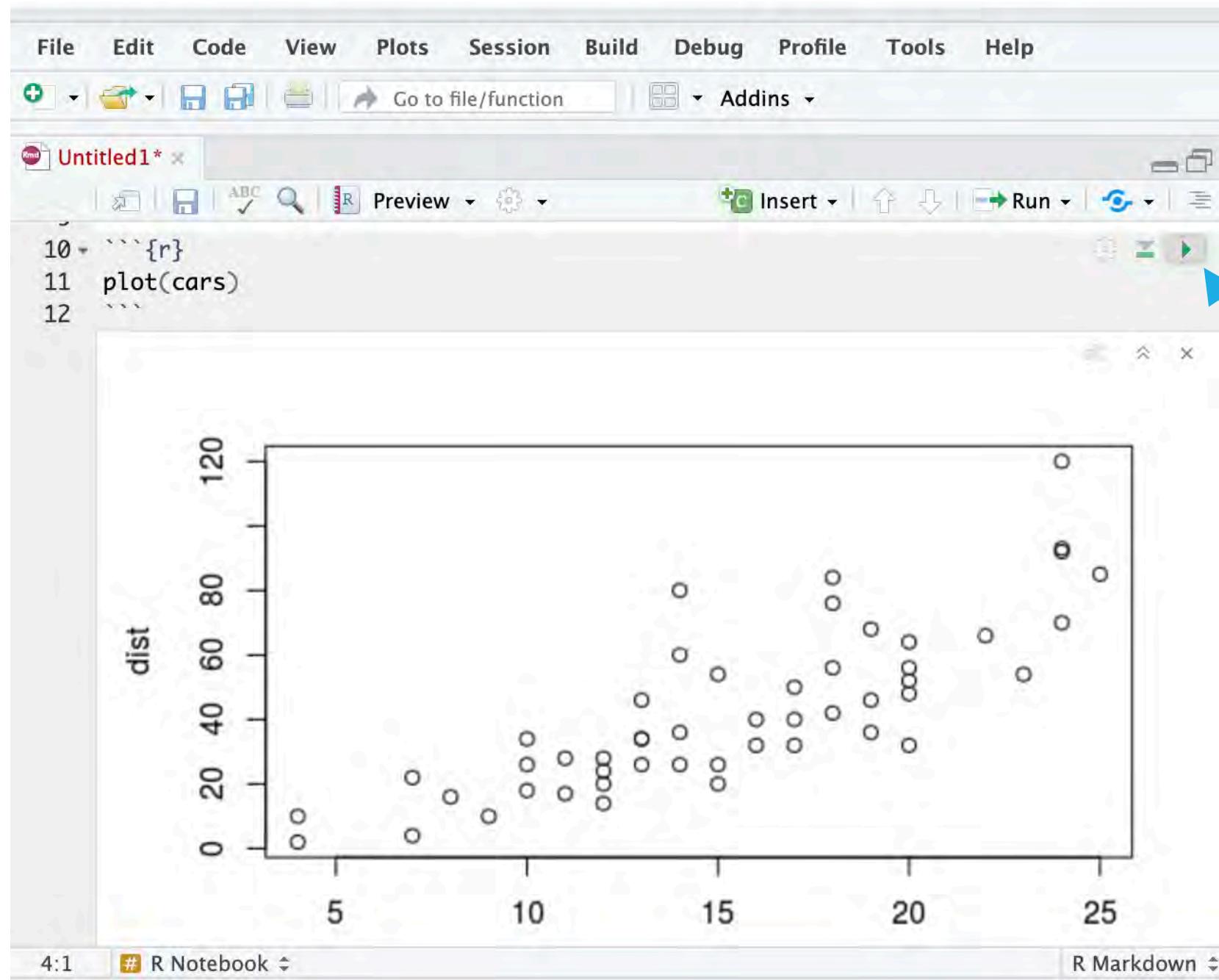
Chunk
name

Run chunk

Code in body of
chunk

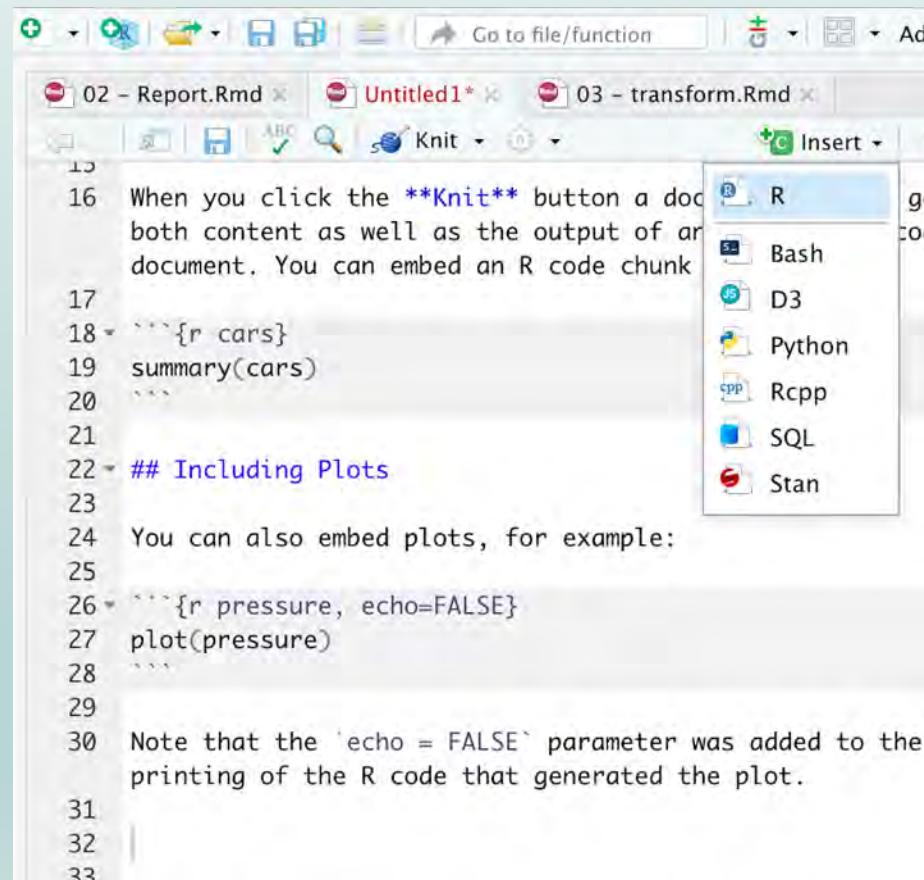
```
17
18  ````{r cars}
19  summary(cars)
20  ```
21
```





Run code chunk

Your Turn #3



The screenshot shows the RStudio interface. The code editor contains the following R code:

```
16 When you click the **Knit** button a doc  
both content as well as the output of an  
document. You can embed an R code chunk  
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.  
31  
32  
33
```

A context menu is open over the code at line 22, showing options like R, Bash, D3, Python, Rcpp, SQL, and Stan.

1. Under your new “My Calculation” heading, insert a code chunk into white space using Insert button on top right of code window
2. Type the following in the grey box of your new code chunk:
`mean(c(10, 20, 30))`
3. Execute code chunk by pressing Run button on top right of code chunk



Rmd Untitled1*

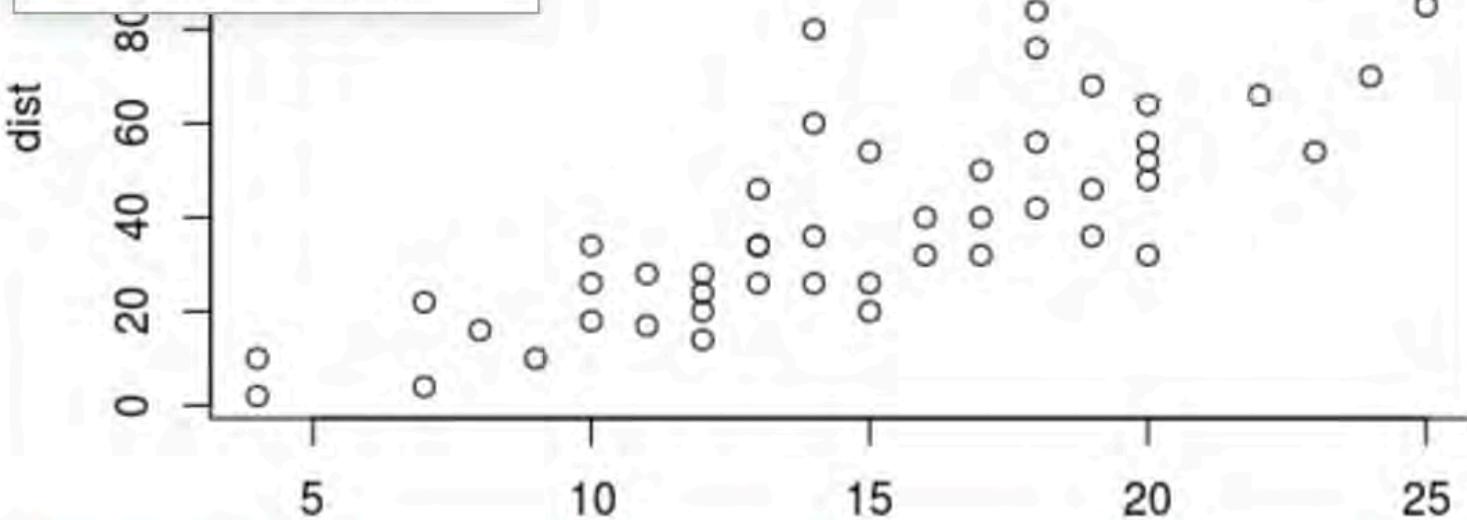
10 +
11 pl
12

- Preview Notebook
- Knit to HTML
- Knit to PDF
- Knit to Word

Knit with Parameters...

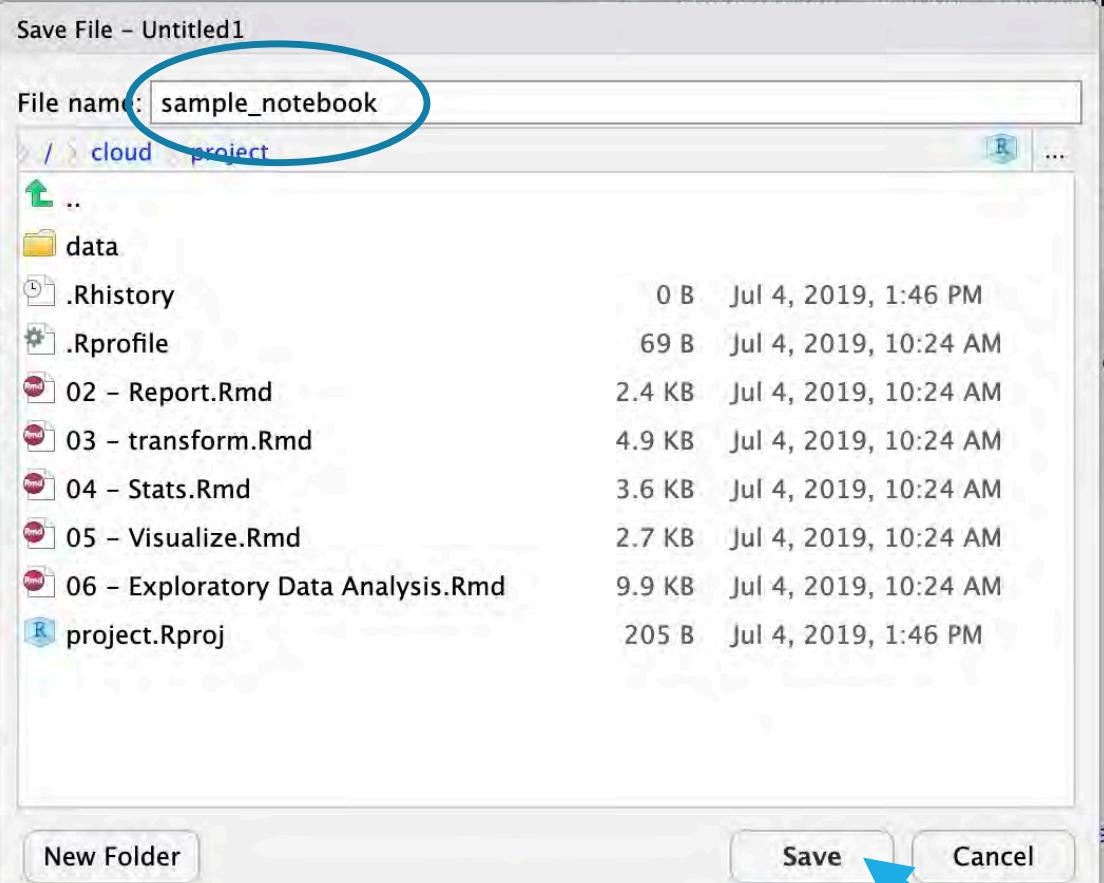
Knit Directory ▾

Clear Knitr Cache...



Untitled1*

```
2 title: "R Notebook"
3 output:
4   html_document:
5     df_print: paged
6 ---
7
8 This is an [R Markdown](http://rmarkdown
code within the notebook, the results ap
9
10 Try executing this chunk by clicking the
placing your cursor inside it and pressi
11
12 ``{r}
13 plot(cars)
14 ``
```



6:1 # R Notebook

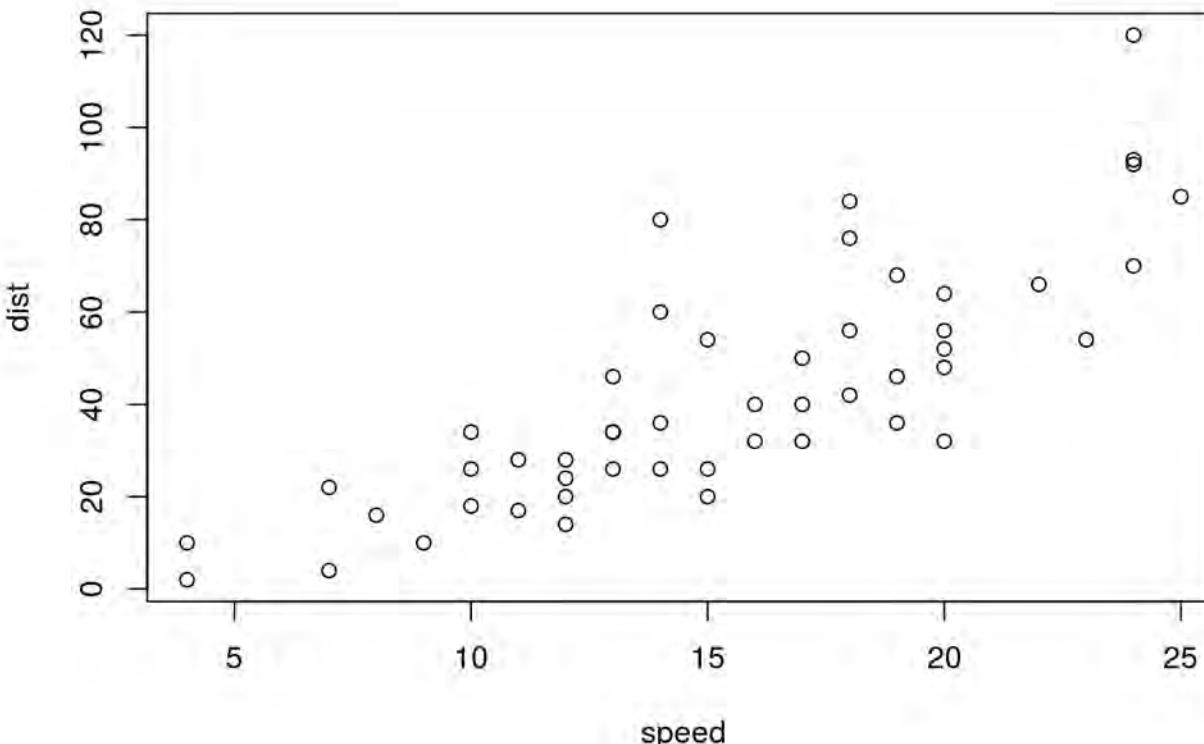
Console

R Notebook

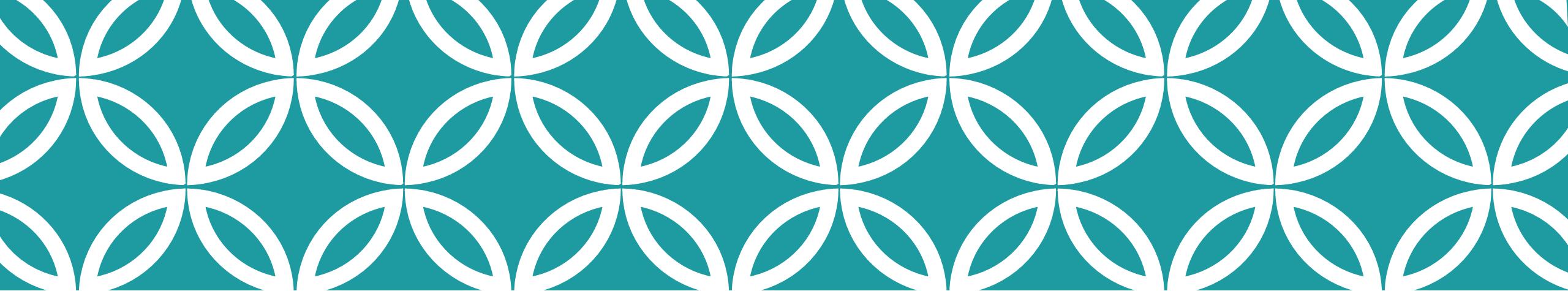
This is an R Markdown Notebook. When you execute code within the notebook, the results appear beneath the code.

Try executing this chunk by clicking the *Run* button within the chunk or by placing your cursor inside it and pressing *Cmd+Shift+Enter*.

```
plot(cars)
```



Add a new chunk by clicking the *Insert Chunk* button on the toolbar or by pressing *Cmd+Option+I*.



What Else?

Why name Code Chunks?

Jump between
chunks

The screenshot shows an RStudio interface with a code editor and a terminal. In the code editor, several lines of R code are visible, including a multi-line comment block. A tooltip is displayed over the code, listing various chunk types: Sample Markdown, R Markdown, Including Plots, and Sample Markdown again. The 'Including Plots' option is highlighted. The terminal below shows the output of the R code, which includes a warning about the 'ALSE' parameter.

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

Console Terminal × Jobs ×
/cloud/project/

R version 3.5.2 (2018-12-20) -- "Froshell Talon"

“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:

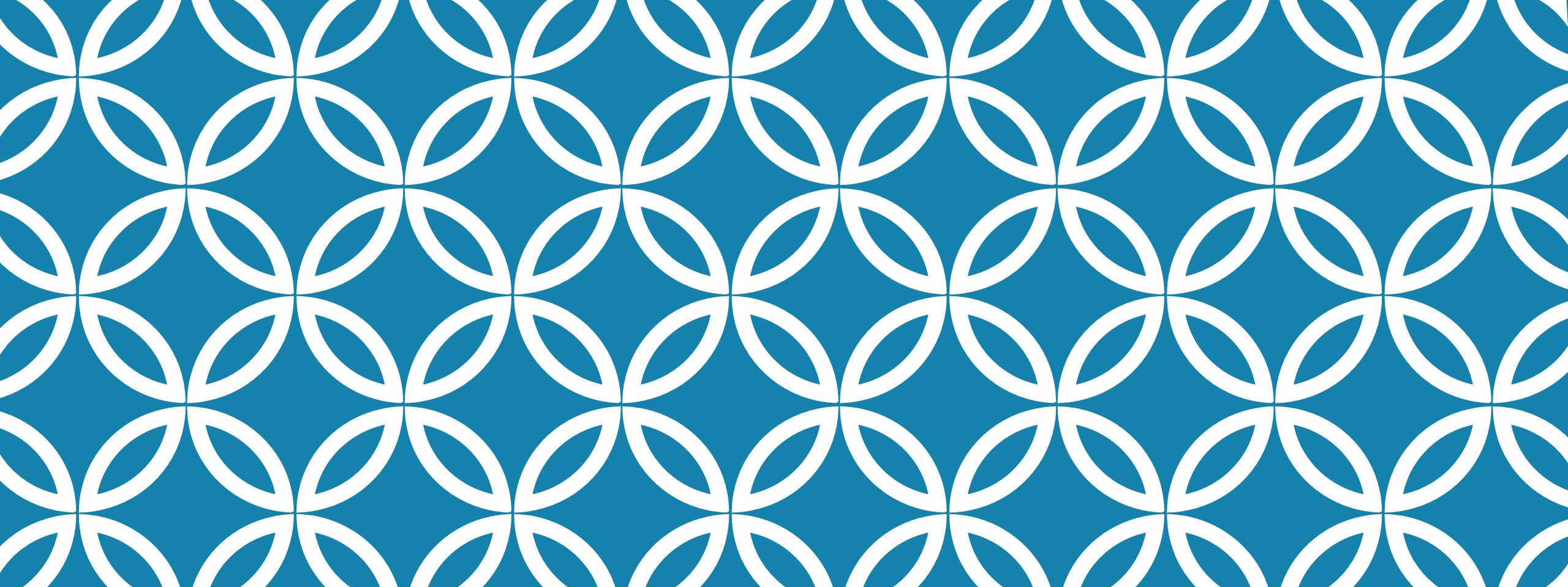
- `install.packages('tinytex')`
- `tinytex::install_tinytex()`

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
Joseph Rudolf

December 13, 2020

Visualization Lecture Courtesy of 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



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

covid_testing

	mrn	first_name	last_name	gender	pan_day	test_id	clinic_name	result
1	5001412	jhezane	westerling	female		4 covid	inpatient ward a	negative
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	negative
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	wynafryd	seaworth	male		9 covid	clinical lab	negative
11	5008115	patrek	frey	male		9 covid	clinical lab	negative
12	5009309	maege	sand	female		9 covid	medical center	negative
13	5008943	myria	rivers	female		9 covid	picu	negative

Showing 1 to 14 of 15,524 entries, 17 total columns

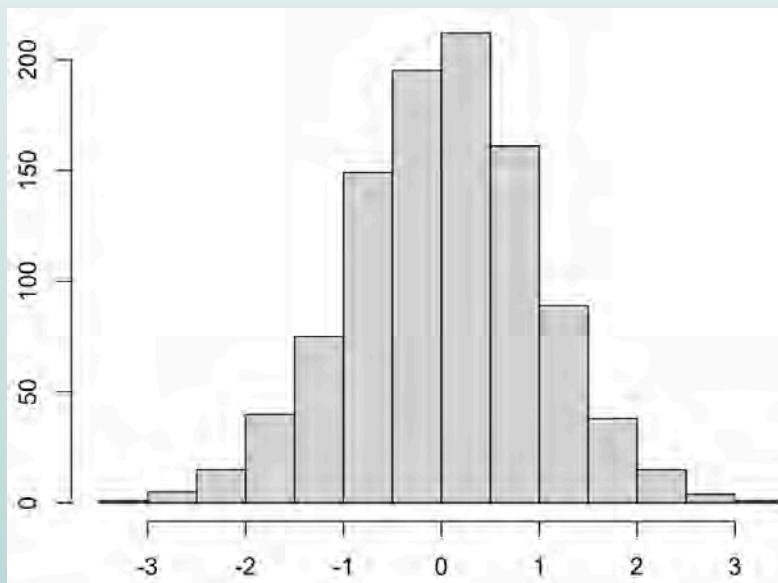
Your Turn 1

Consider the `covid_testing` data frame.

What do you think the 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!

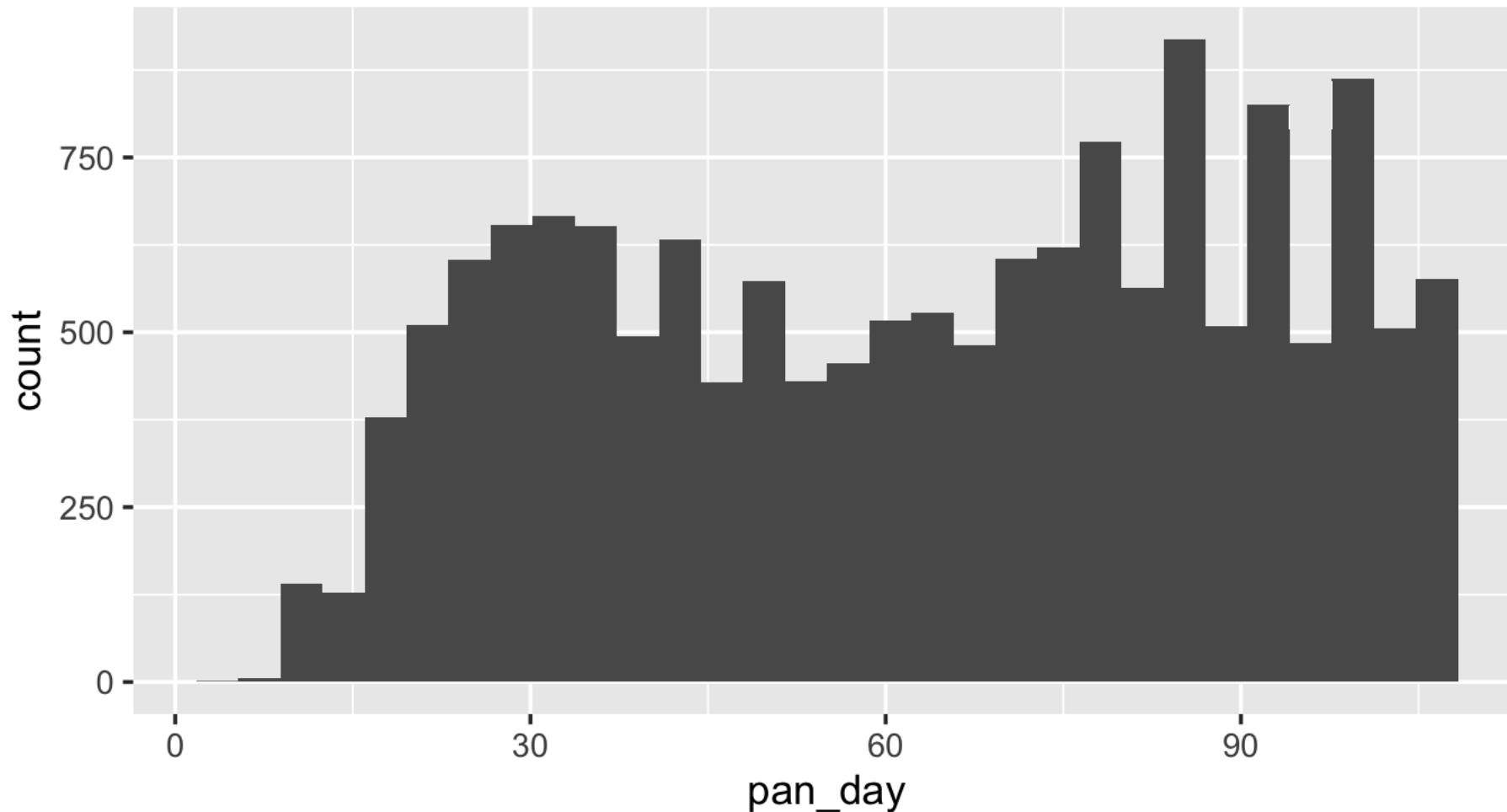
```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day))
```

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

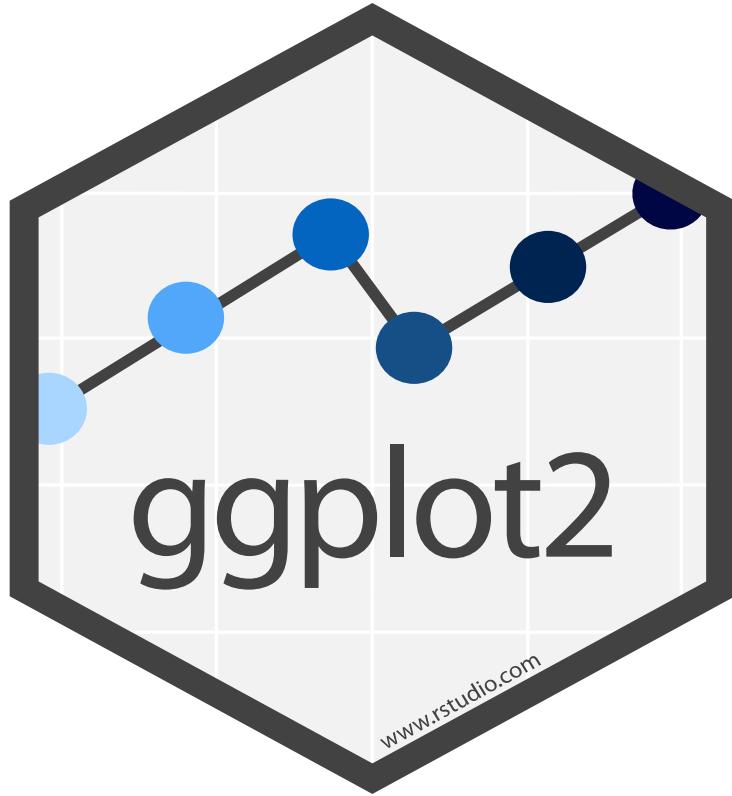
R lets you know that when you ask it to draw a histogram you should tell it how wide each bin should be, because this affects the granularity of the data displayed.

`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

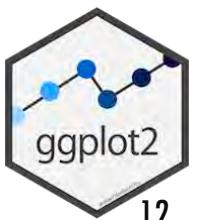
x axis
mapping



To make **any** kind of graph:

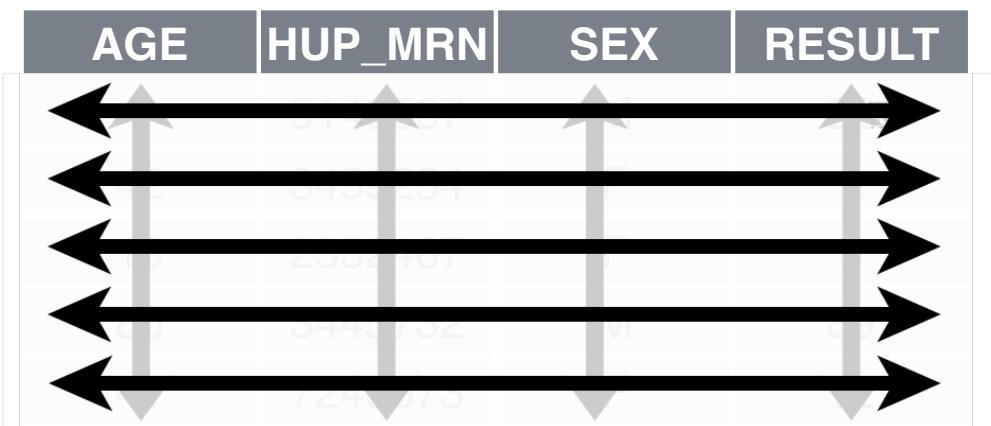
1. Choose a “tidy”
data frame

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



1. Pick a “Tidy” Data Frame

AGE	HUP_MRN	SEX	RESULT
1			
2			
3			
4			
5			



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:

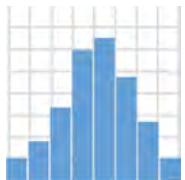
1. Choose a “tidy”
data frame

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

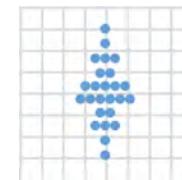
2. Pick a “**geom**”
function



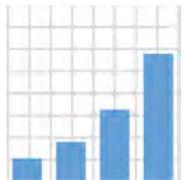
2. Choose a “Geom” Function



`geom_histogram()`



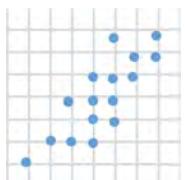
`geom_dotplot()`



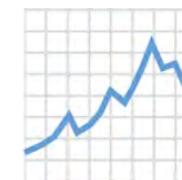
`geom_bar()`



`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>
```

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**

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, ymin, ymax, alpha, color, group, linetype, size**

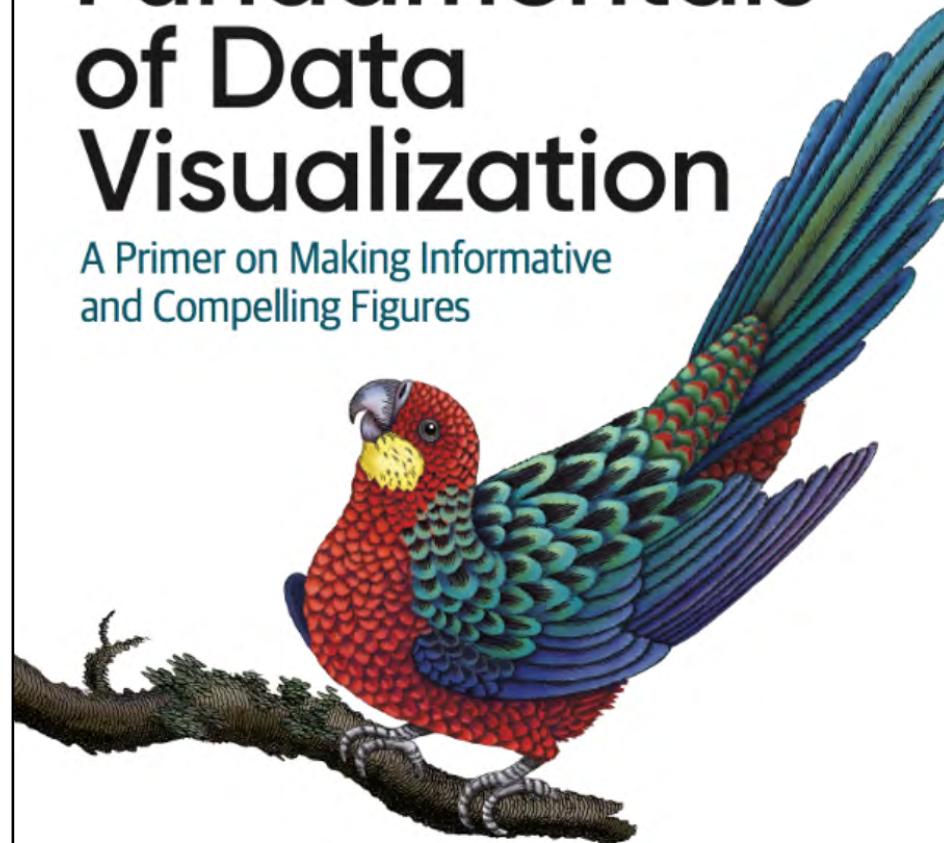
**j + geom_pointrange()
x, y, ymin, ymax, alpha, color, fill, group, linetype,
shape, size**

maps

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

Fundamentals of Data Visualization

A Primer on Making Informative
and Compelling Figures



Claus O. Wilke

<https://serialmentor.com/dataviz>

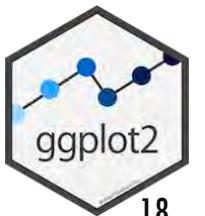
To make **any** kind of graph:

1. Choose a “tidy”
data frame

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

2. Pick a “**geom**”
function

3. Write aesthetic
mappings



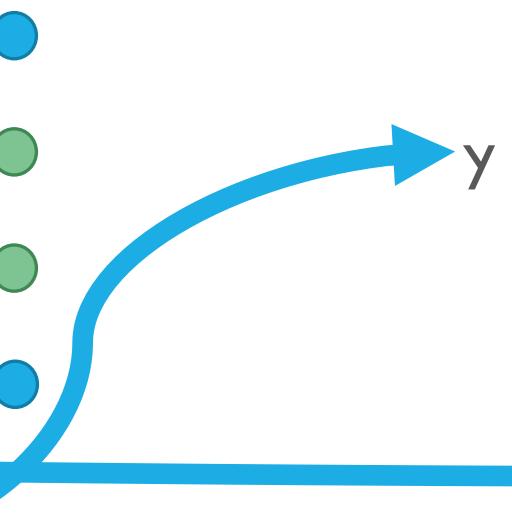
3. Write Aesthetic Mappings

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

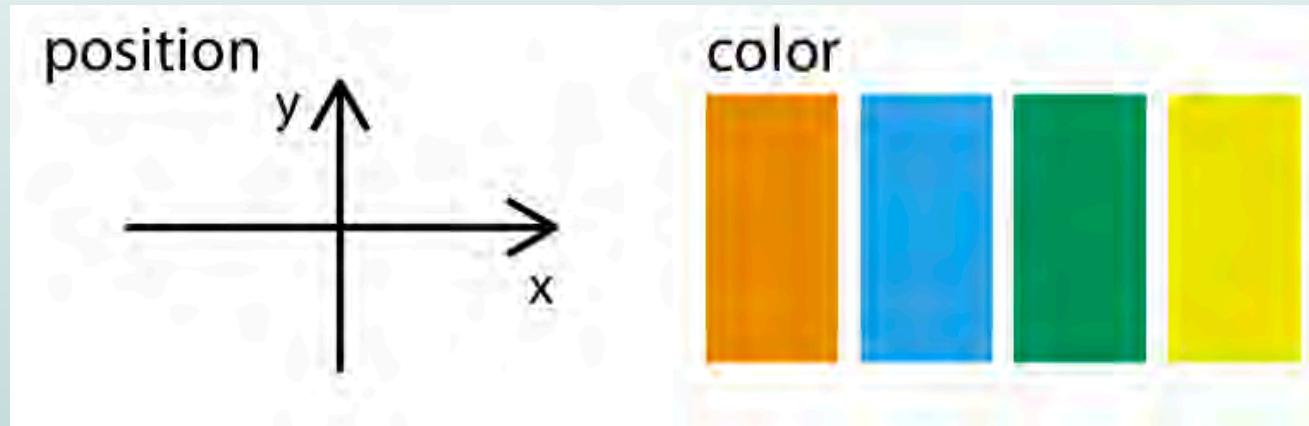
Data frame

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

Graph

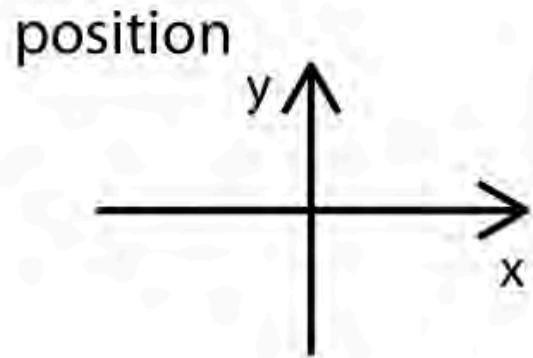


Your Turn 4

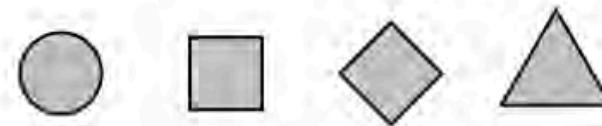


In addition to x/y position and color, what other aesthetics can you think of?
Type your answers in the chat!

Aesthetics



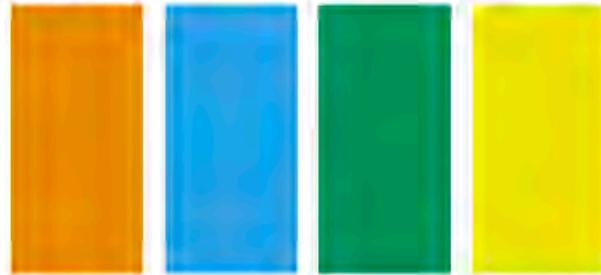
shape



size



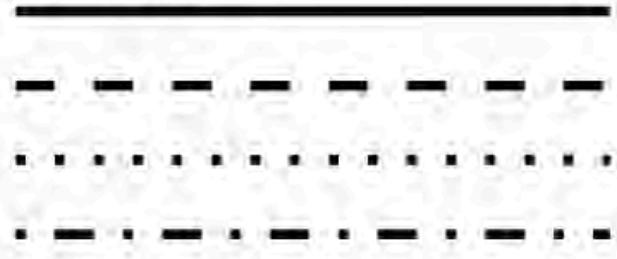
color



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

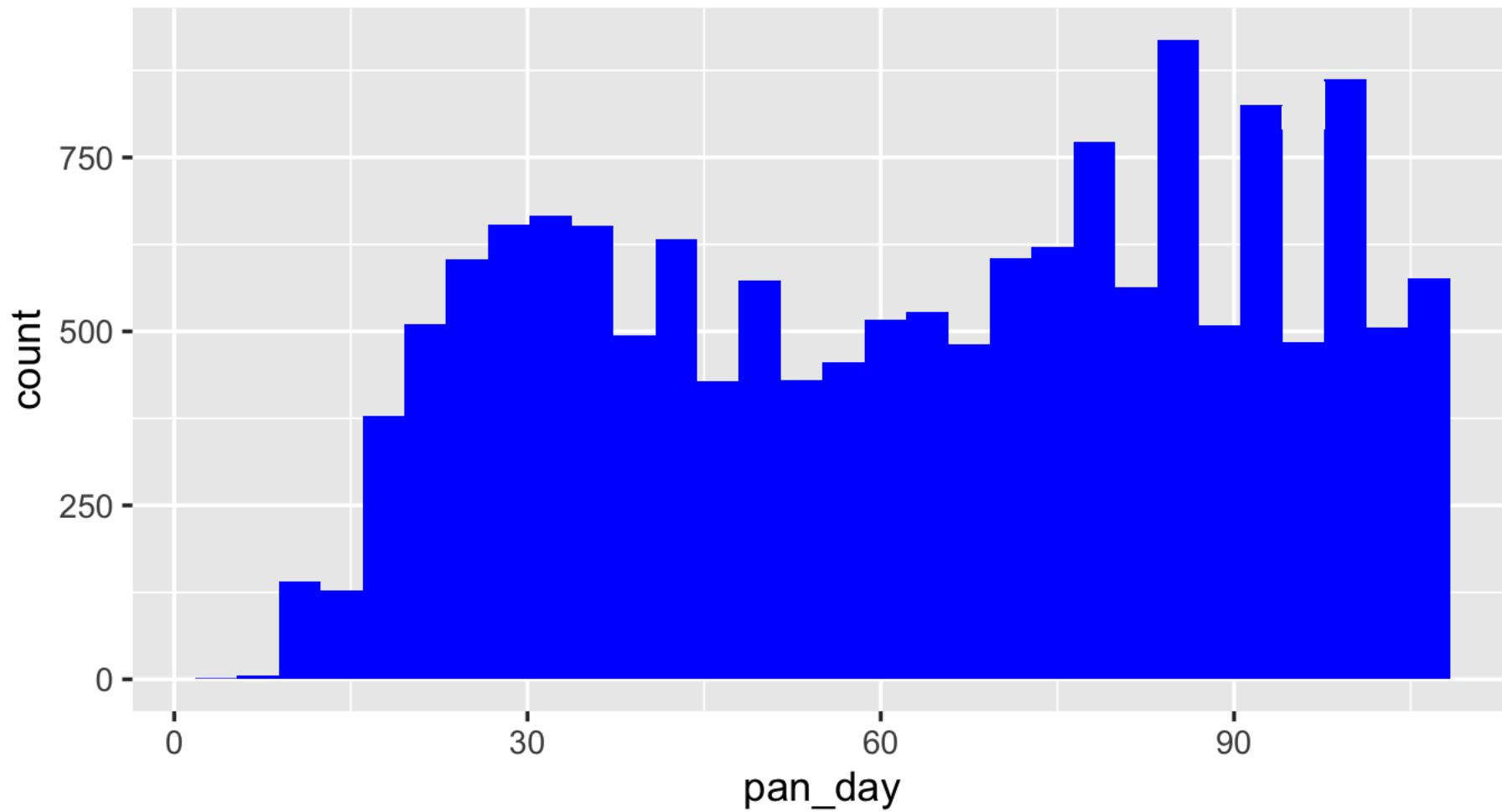


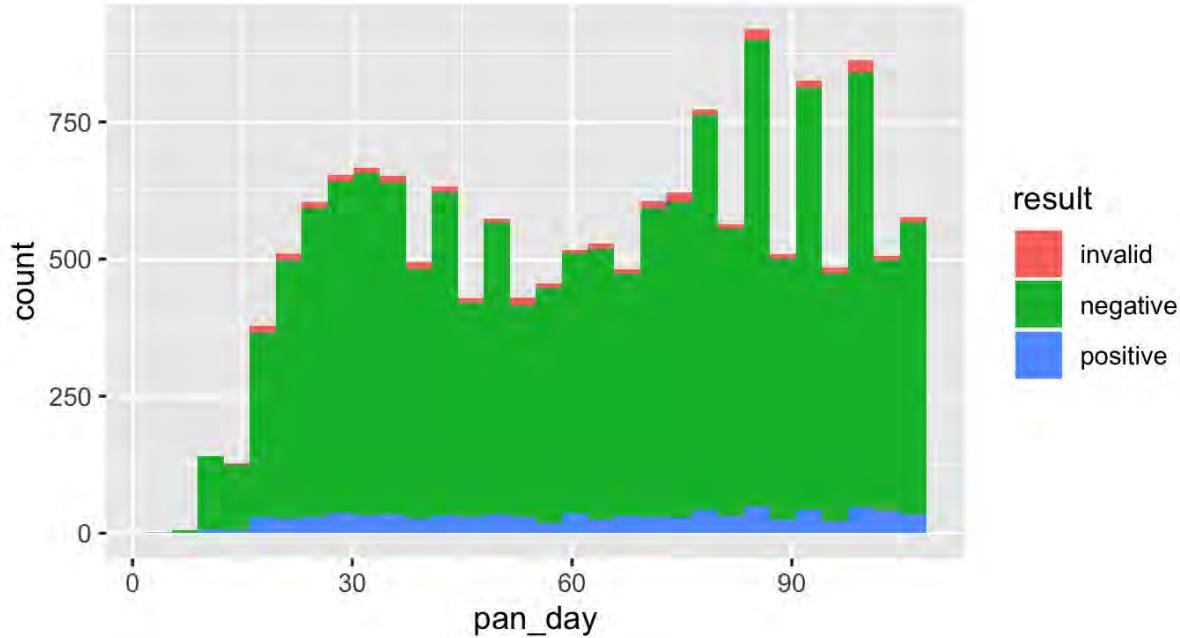
Your Turn 5

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?

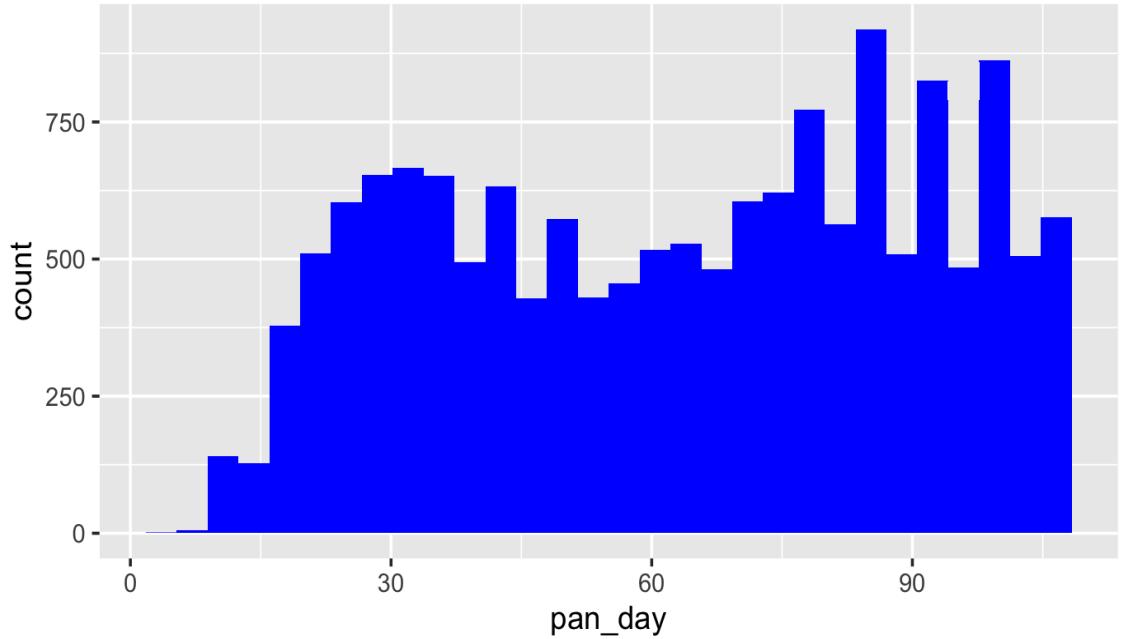




Inside of `aes()`:
map an aesthetic
to a variable

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

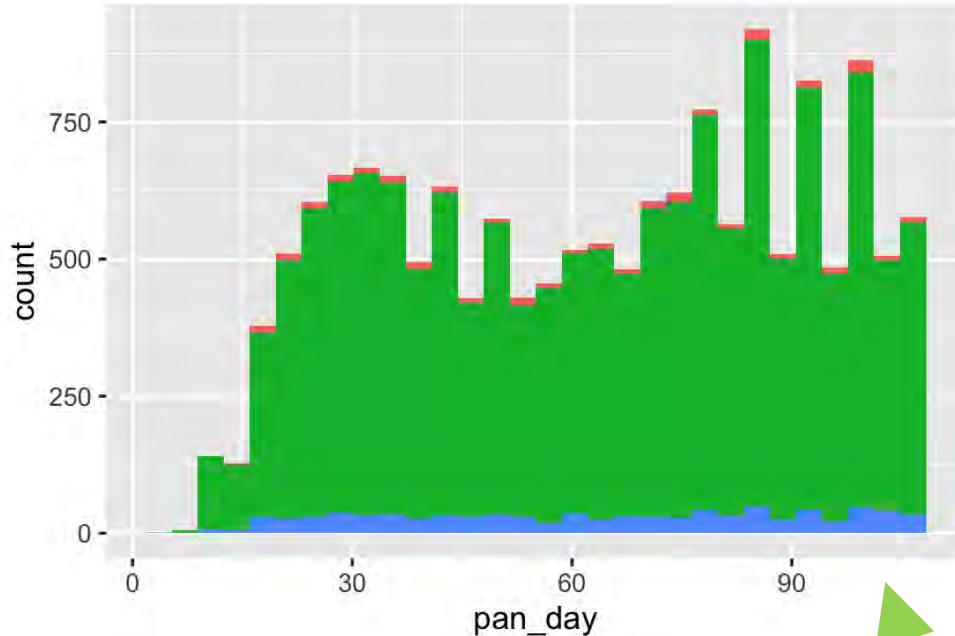




Outside of aes():
set an aesthetic to
a **value**

color name in
“quotes”

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

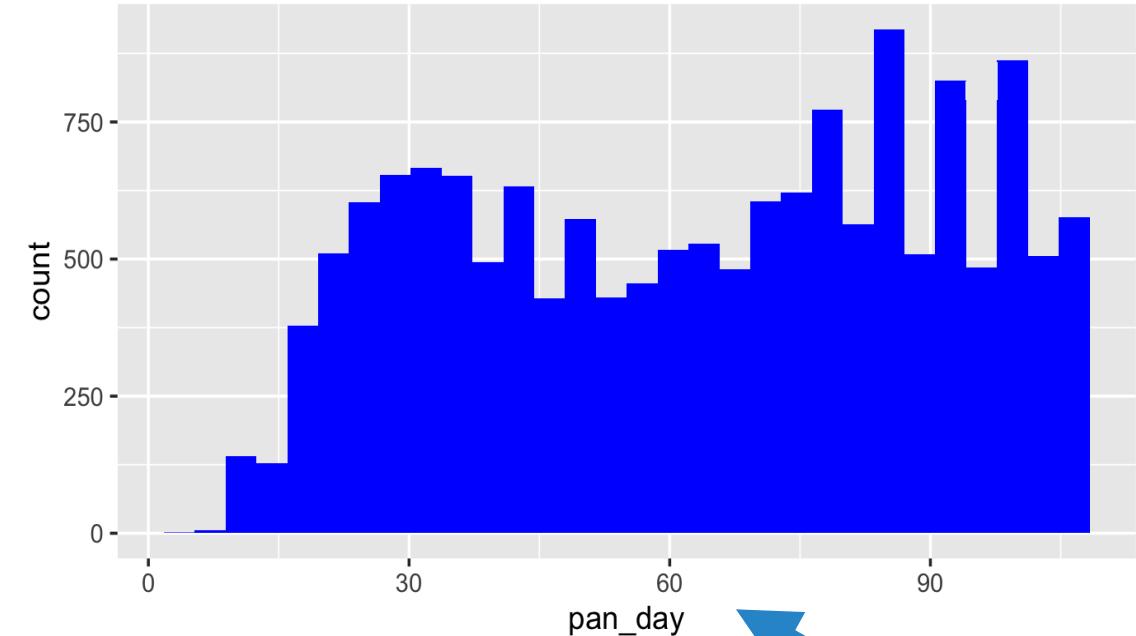


result

invalid

negative

positive

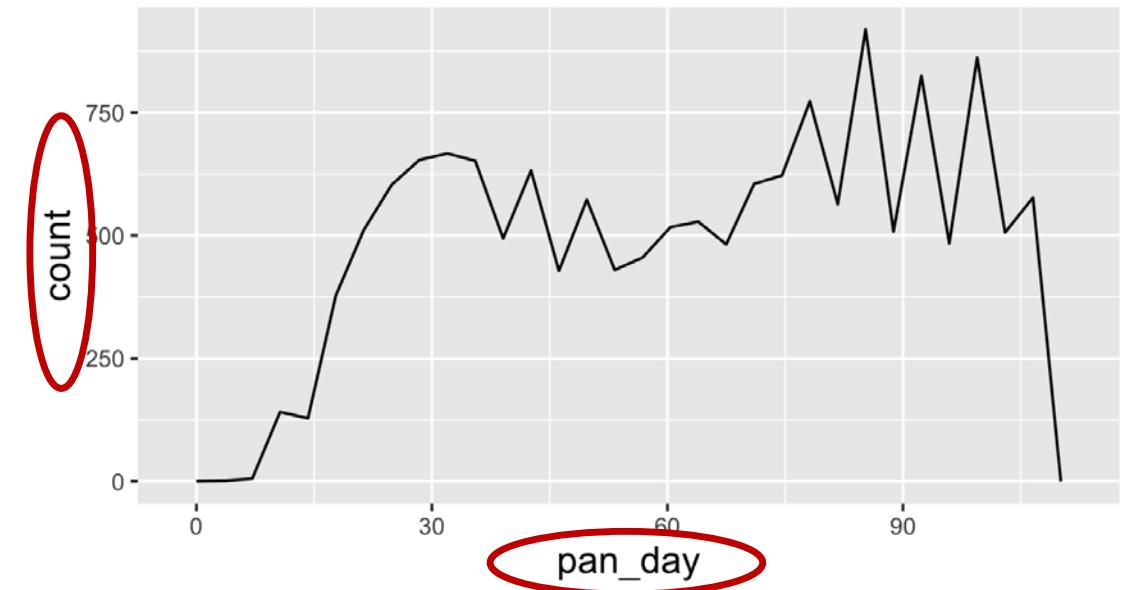
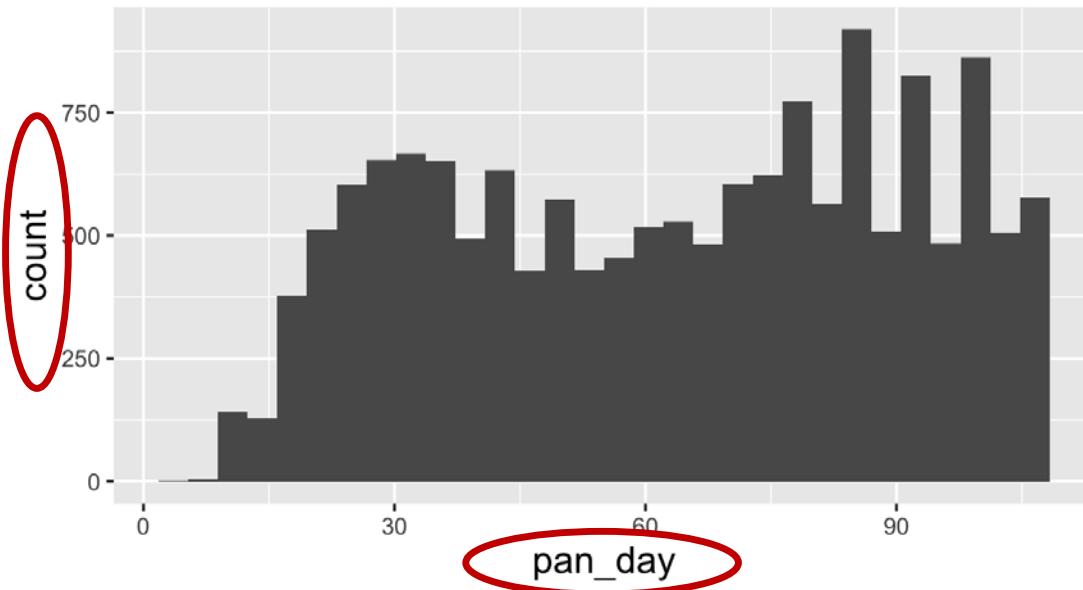


```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day, fill = result))
```

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

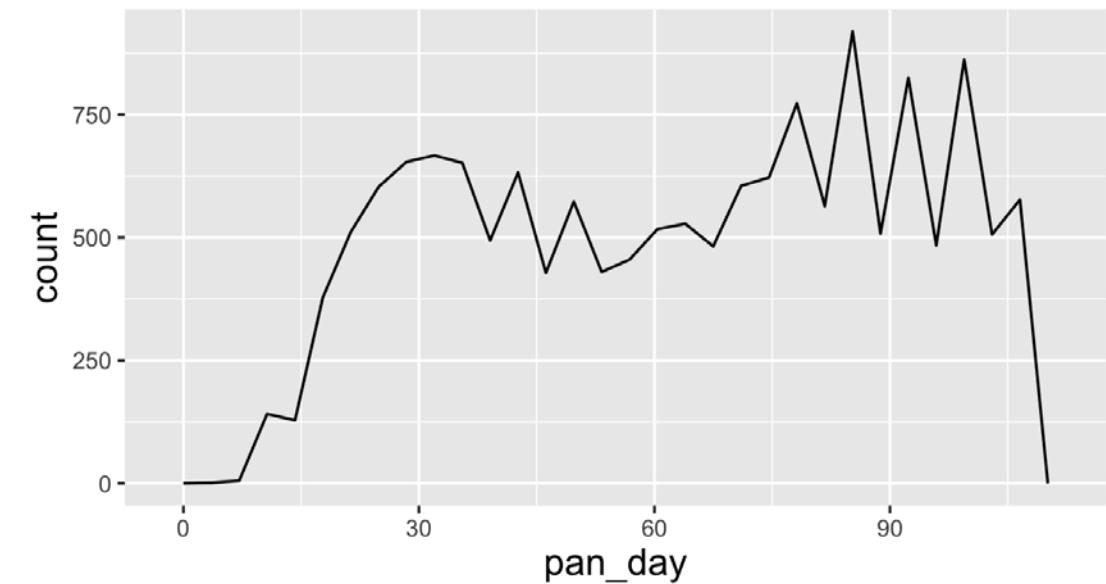
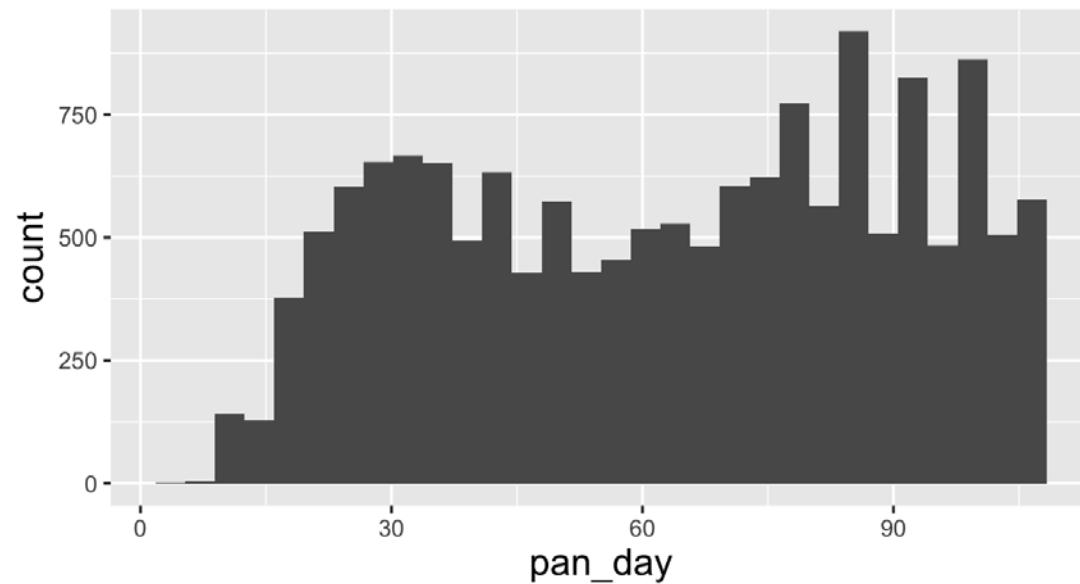
Geom Functions

How are these plots similar?



Same: x axis, y axis, data

How are these plots different?

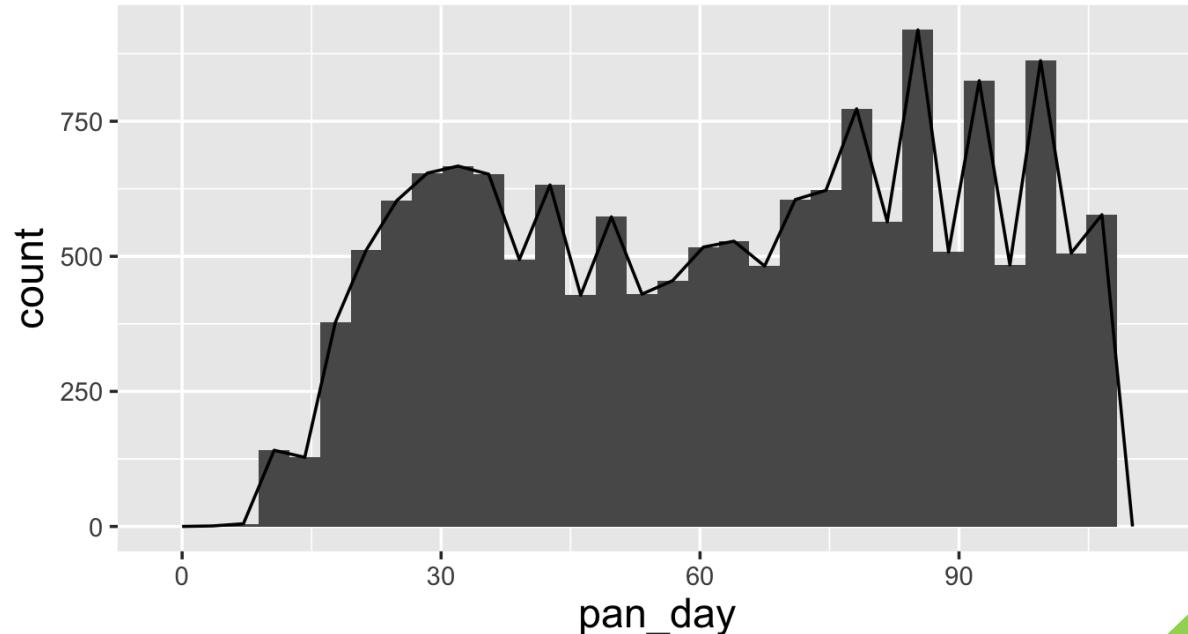


Different **geometric object** (“geom”) used to represent the data

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

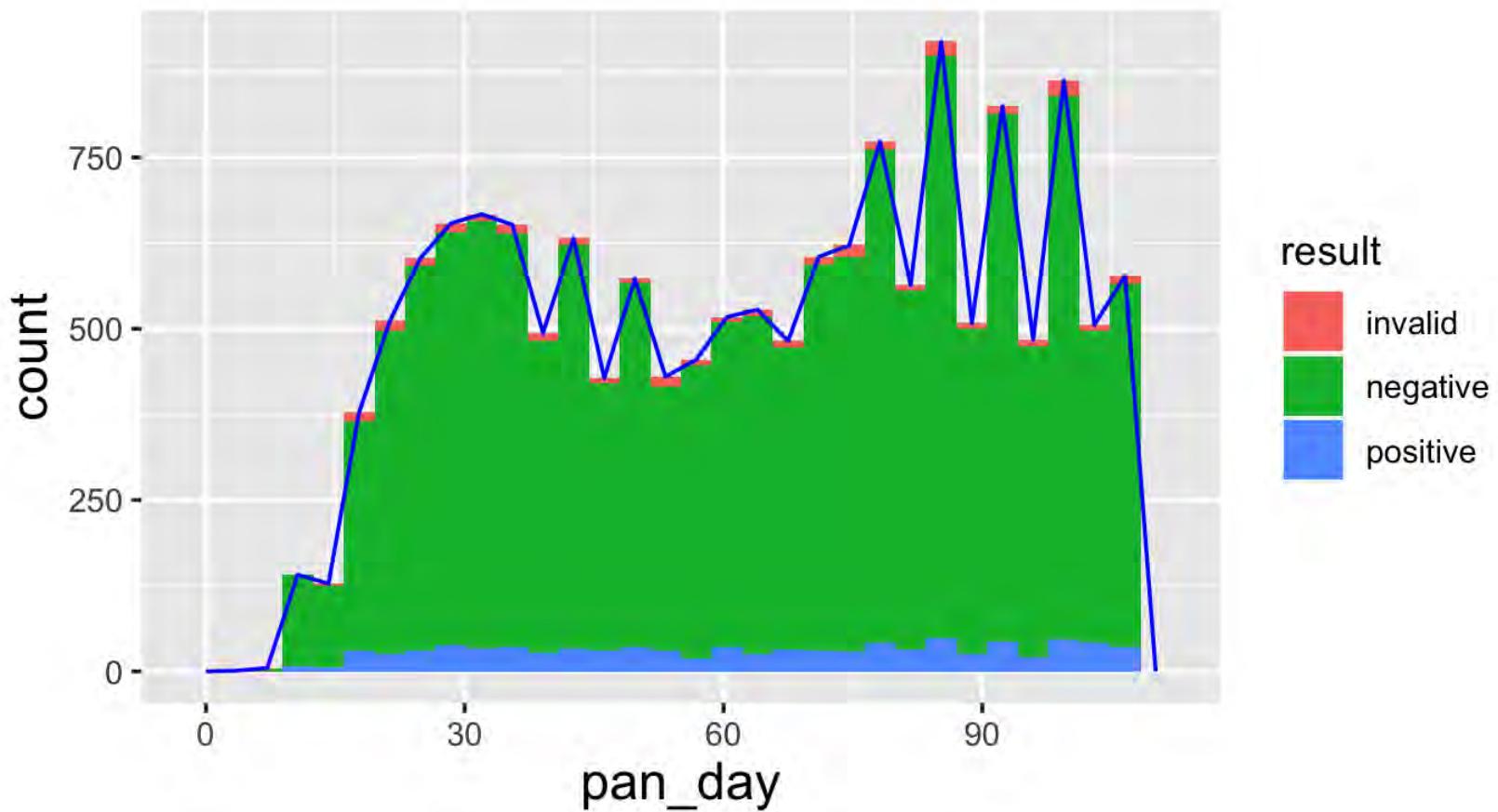
Inside of `ggplot()`:
apply **globally** to
every layer

```
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()
```

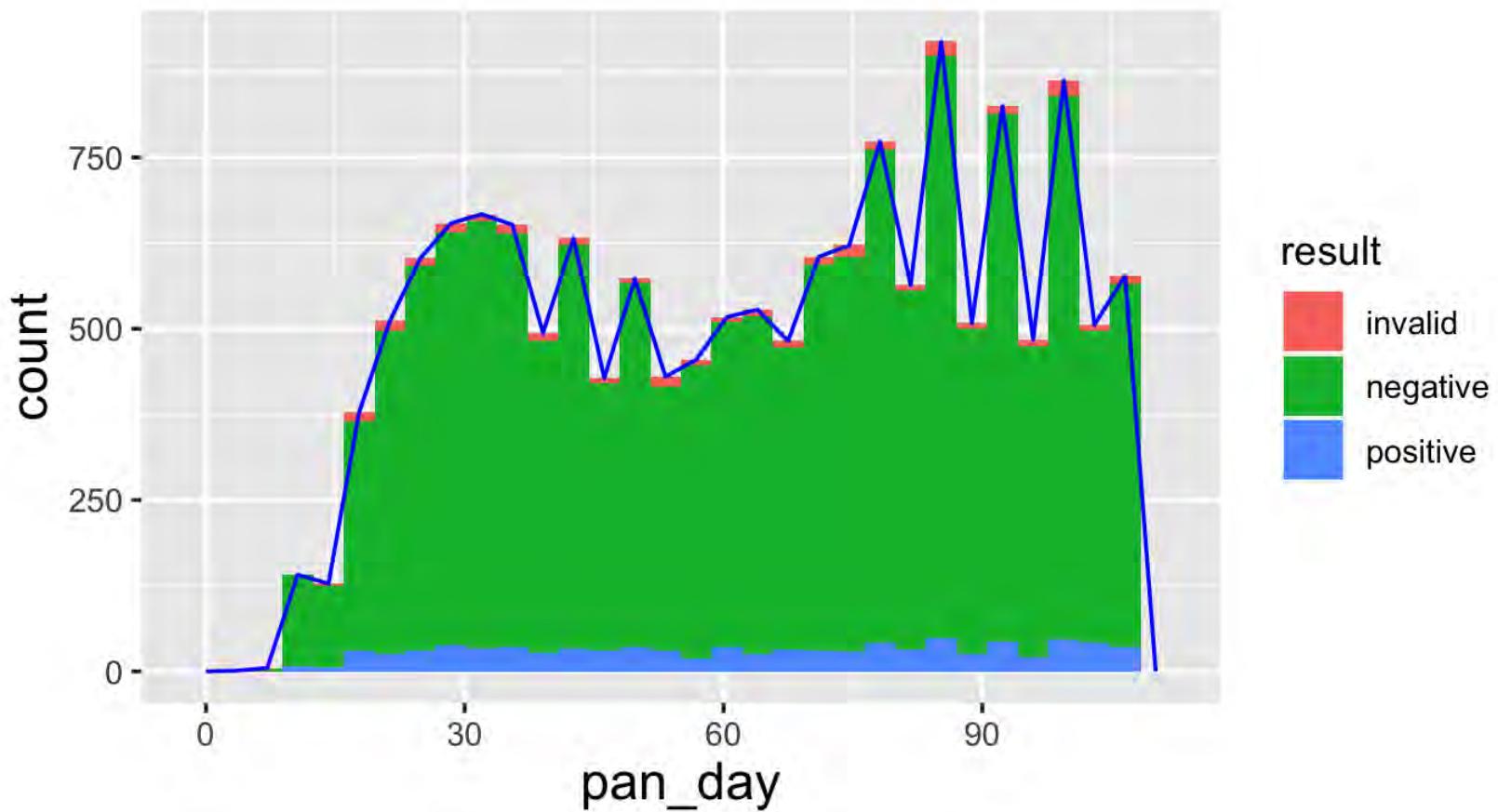
empty ()

How would you make this plot?



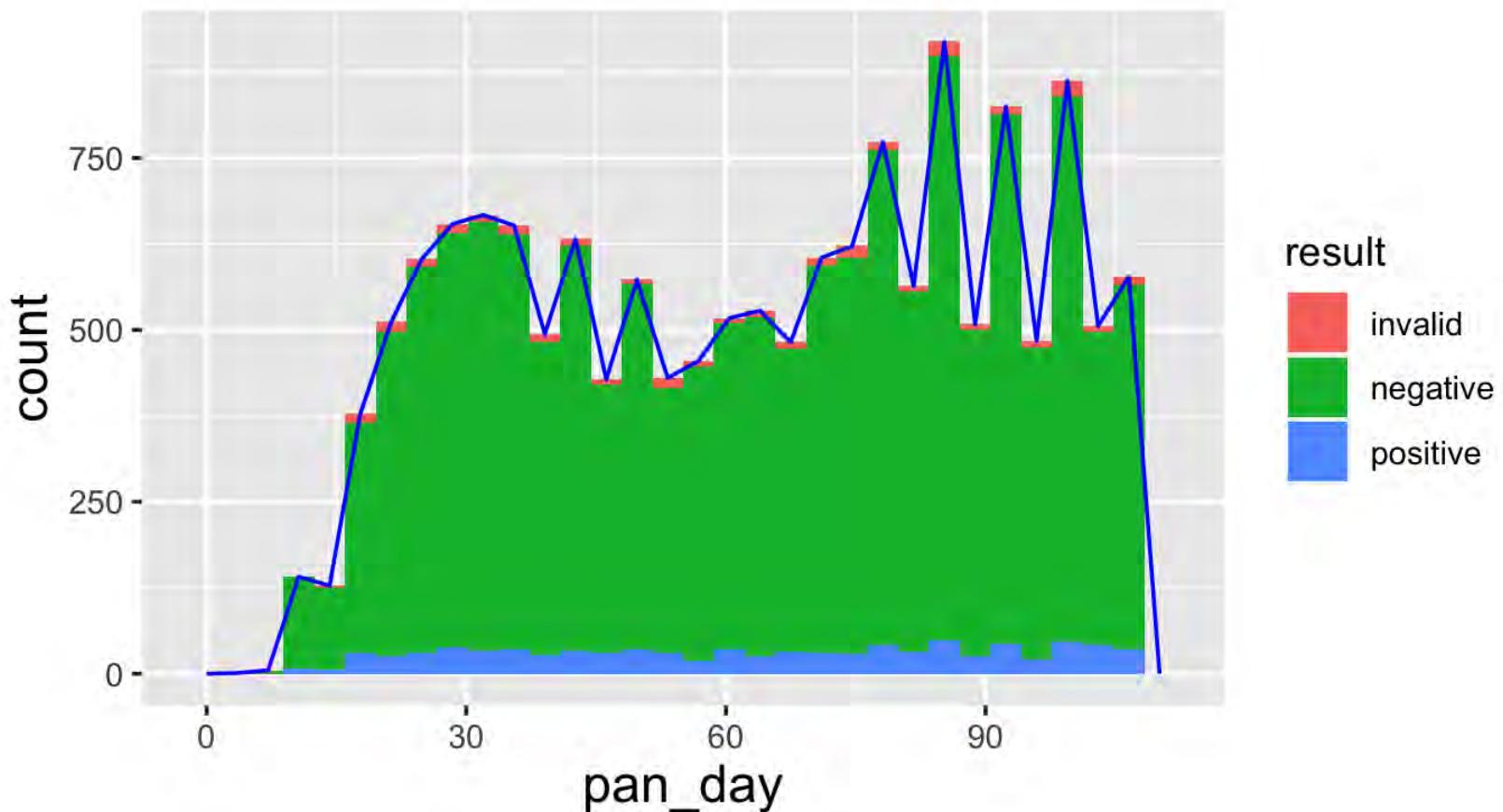
```
ggplot(
```

How would you make this plot?



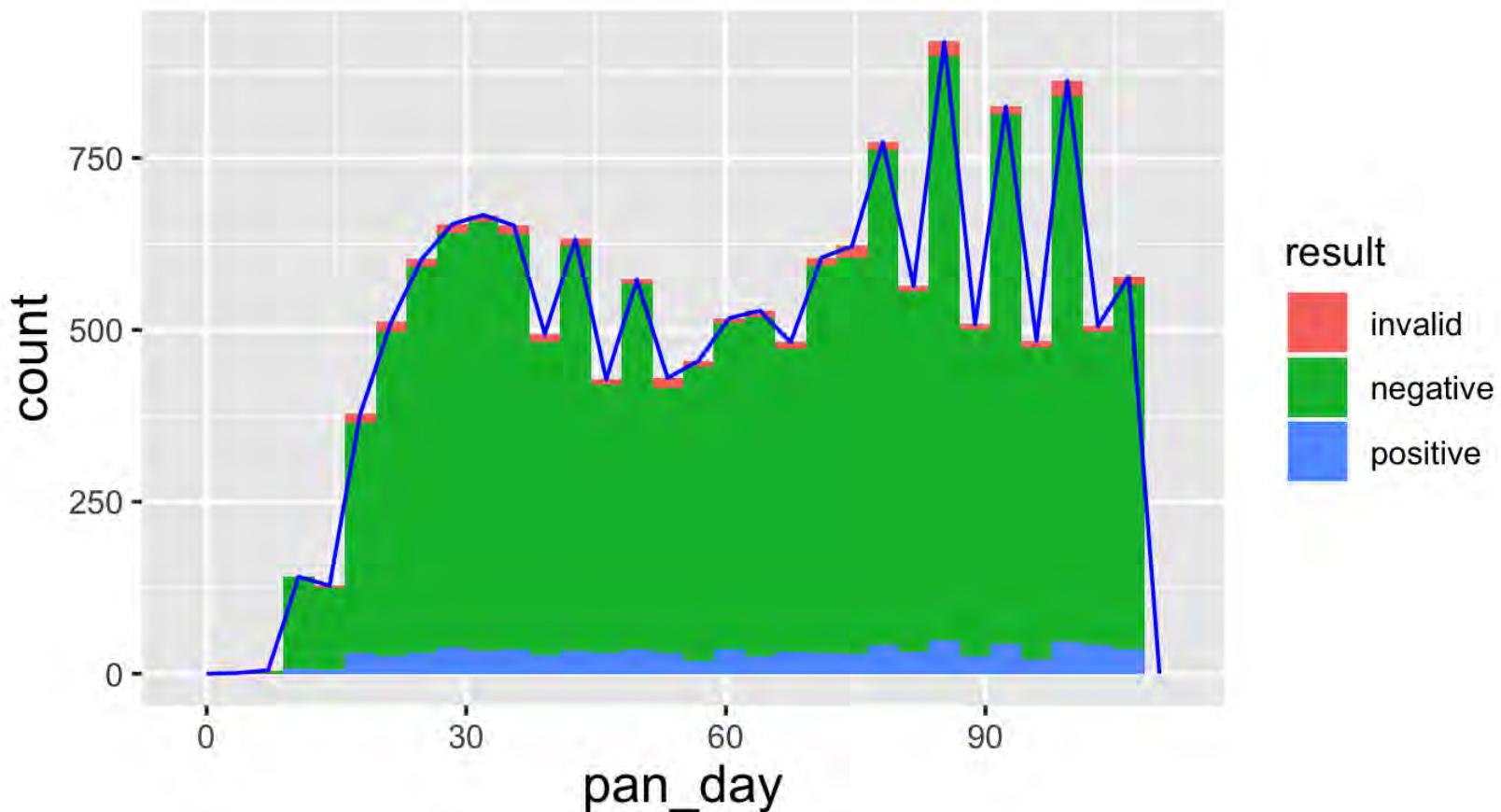
```
ggplot(data = covid_testing
```

How would you make this plot?



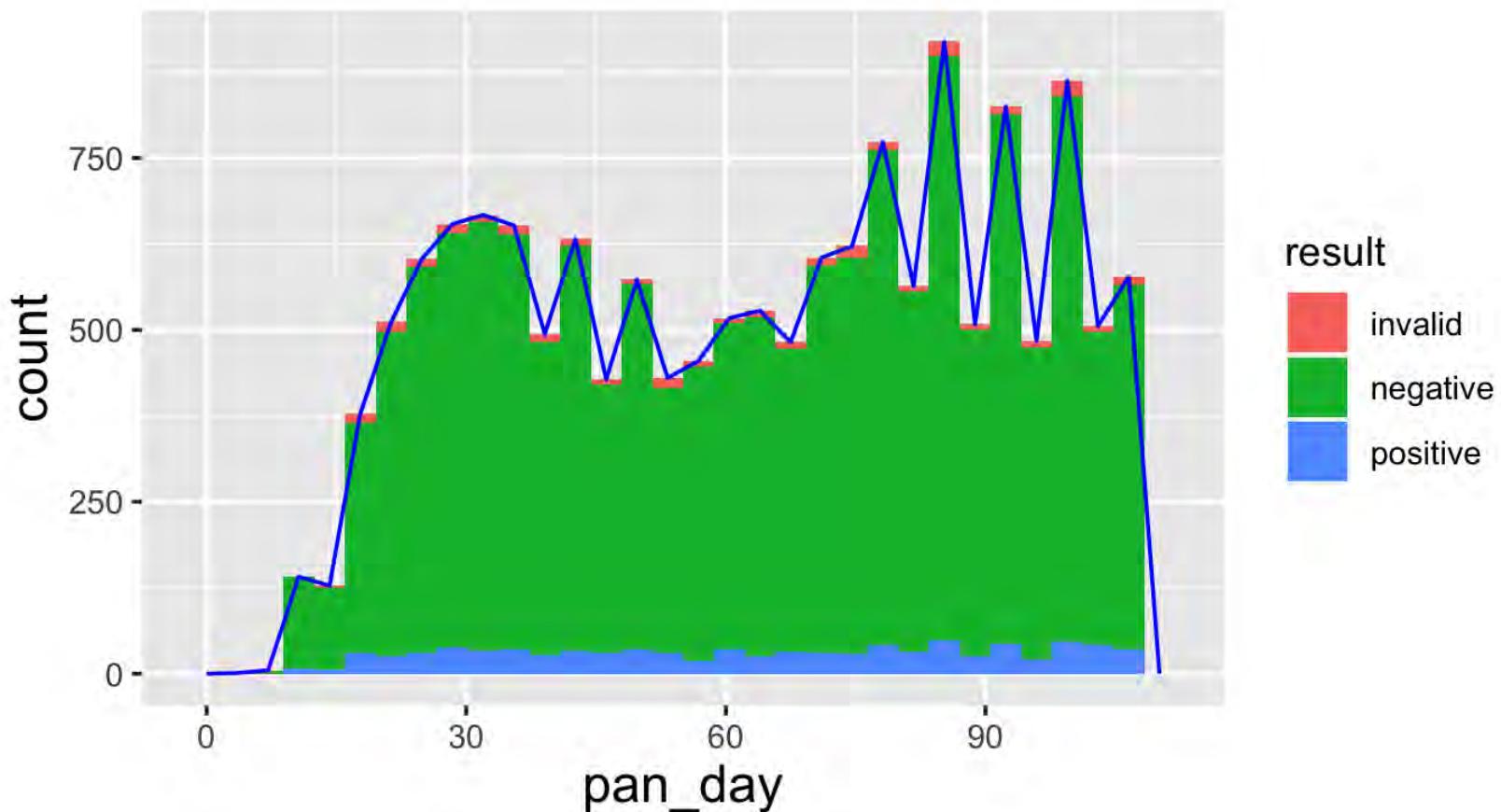
```
ggplot(data = covid_testing, mapping = aes(x = pan_day)) +  
  geom_histogram(
```

How would you make this plot?



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

How would you make this plot?

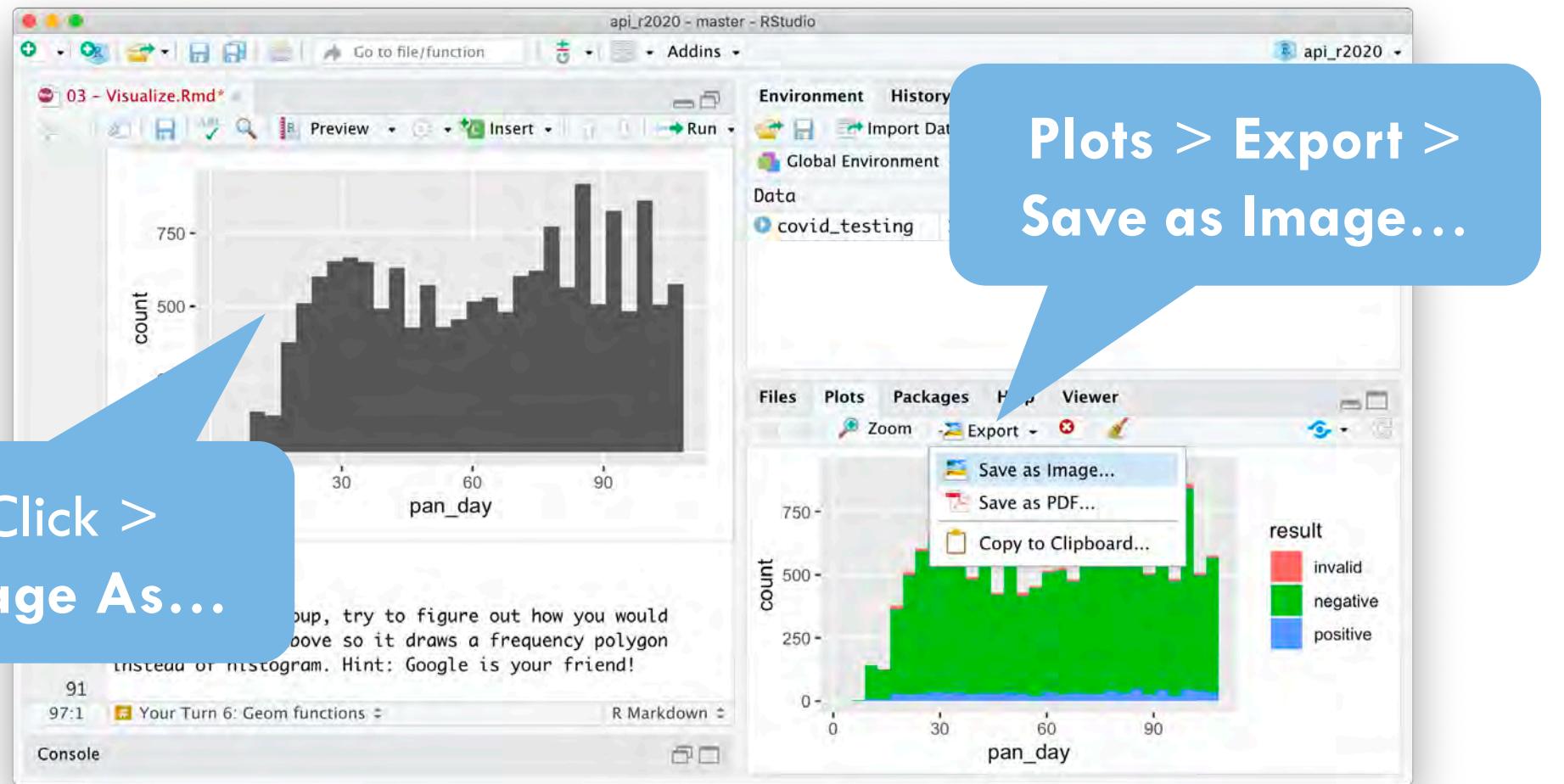


```
ggplot(data = covid_testing, mapping = aes(x = pan_day)) +  
  geom_histogram(mapping = aes(fill = result)) +  
  geom_freqpoly(color = "blue")
```



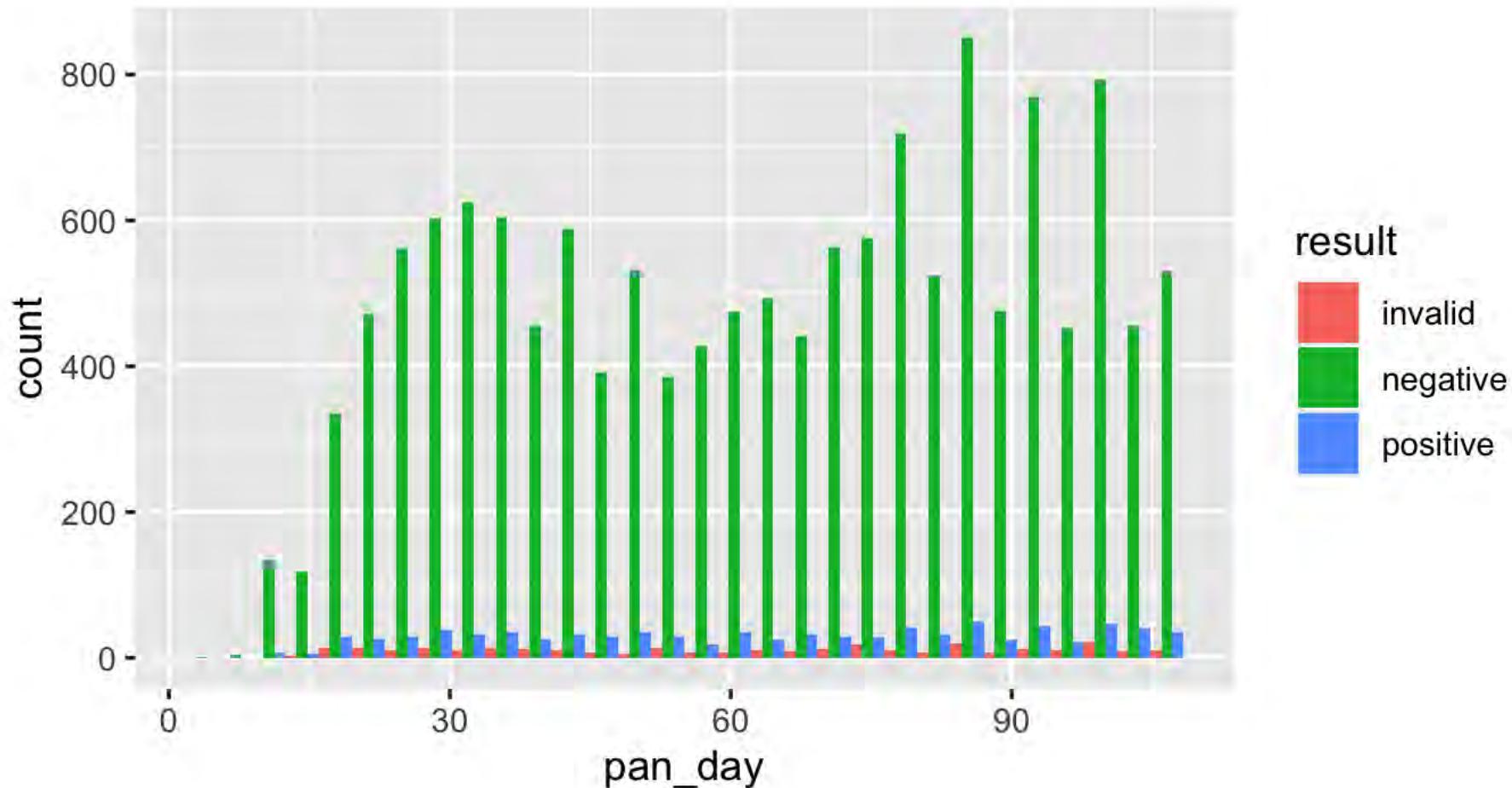
What Else?

Manually saving plots



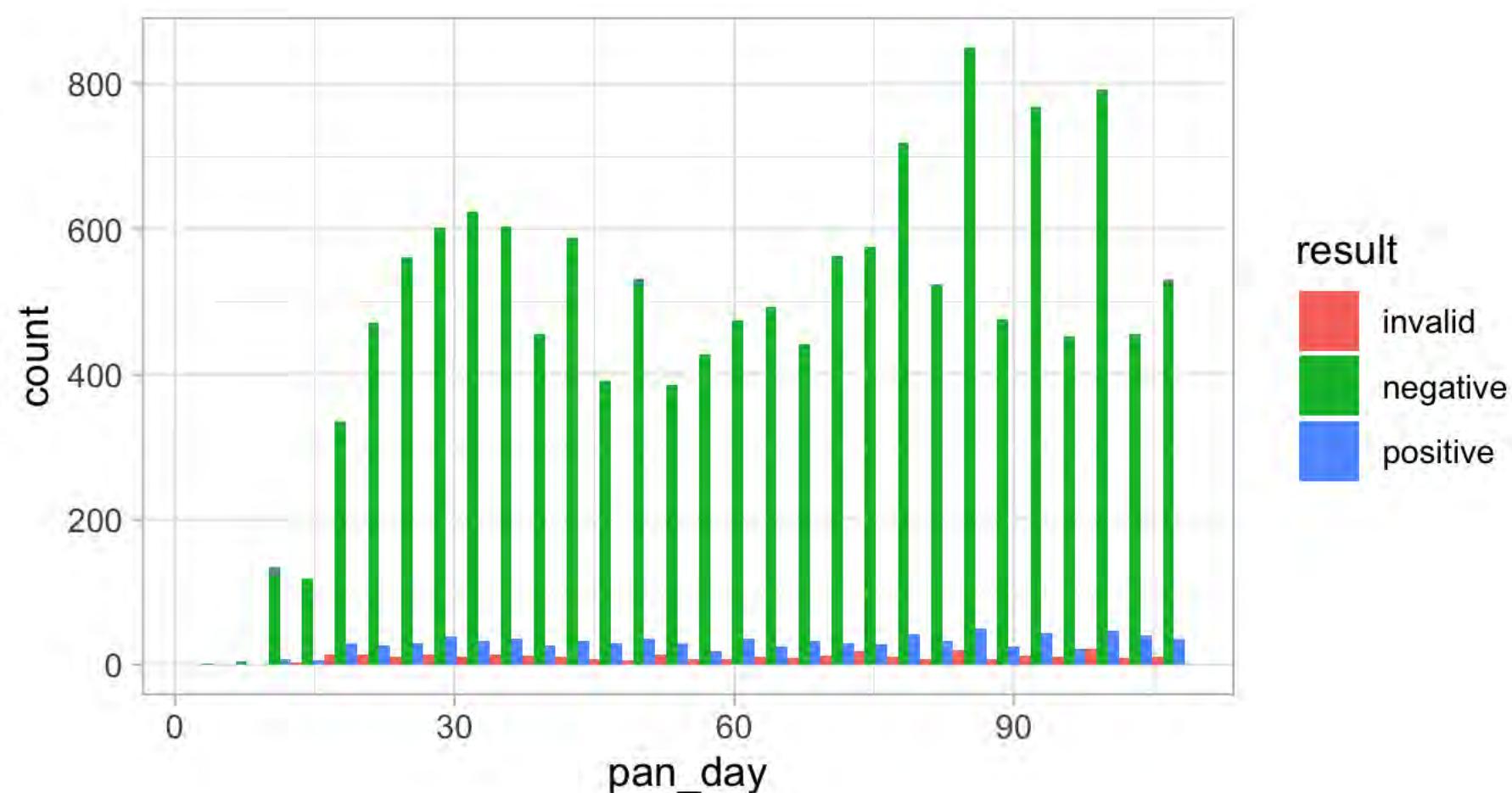
Position adjustments

How overlapping objects are arranged



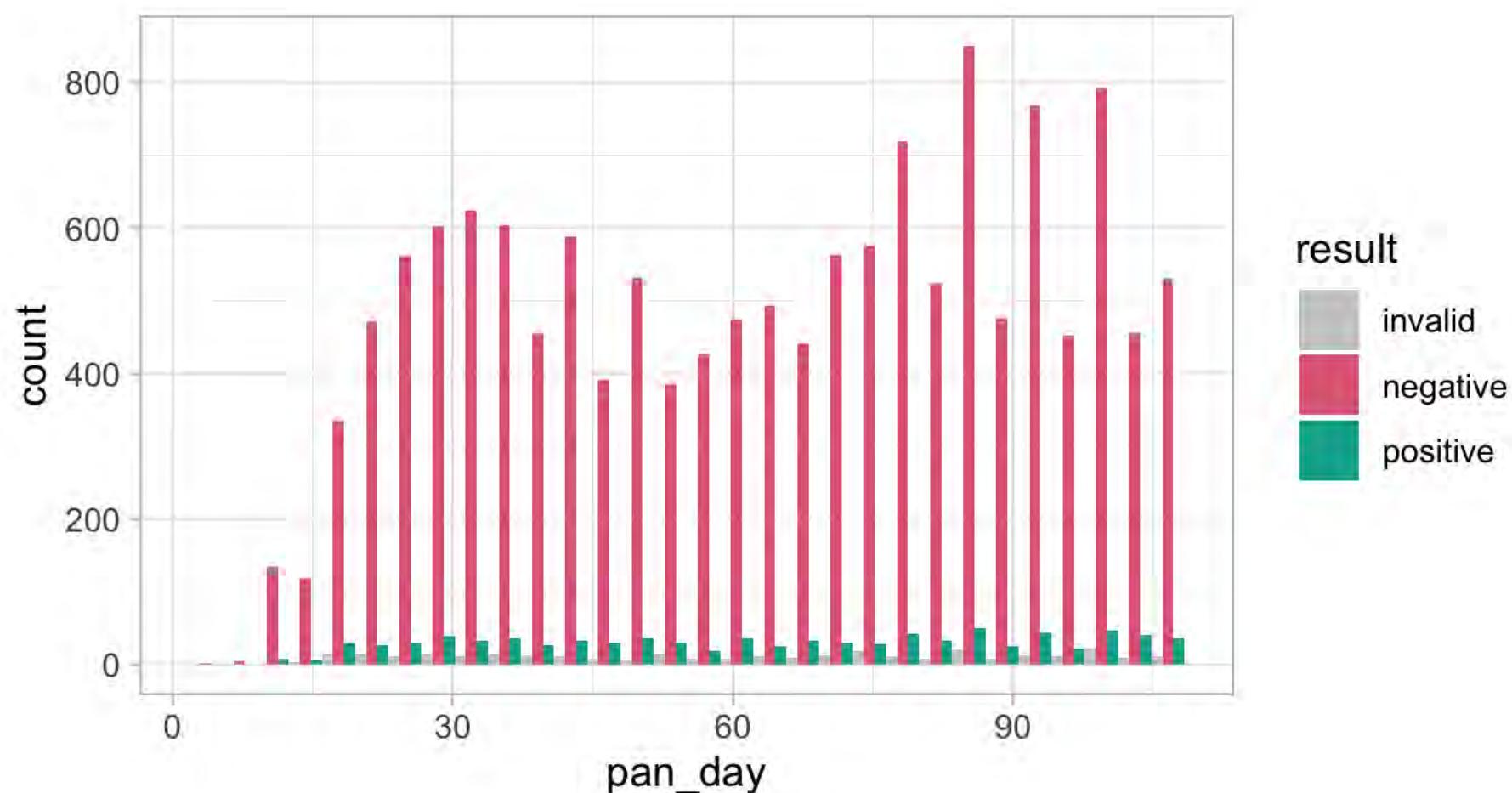
Themes

Visual appearance of non-data elements



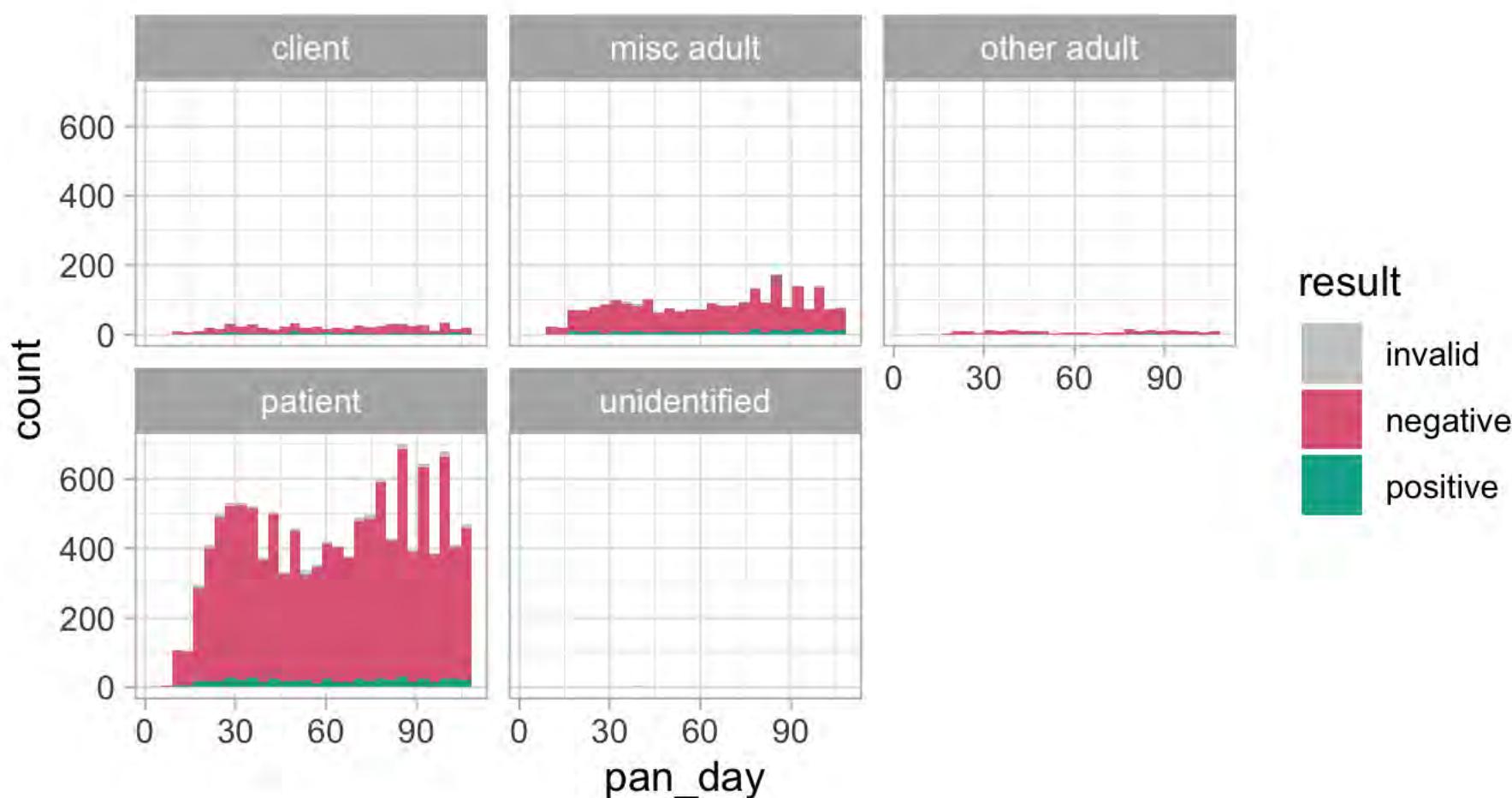
Scales

Customize color scales and other mappings

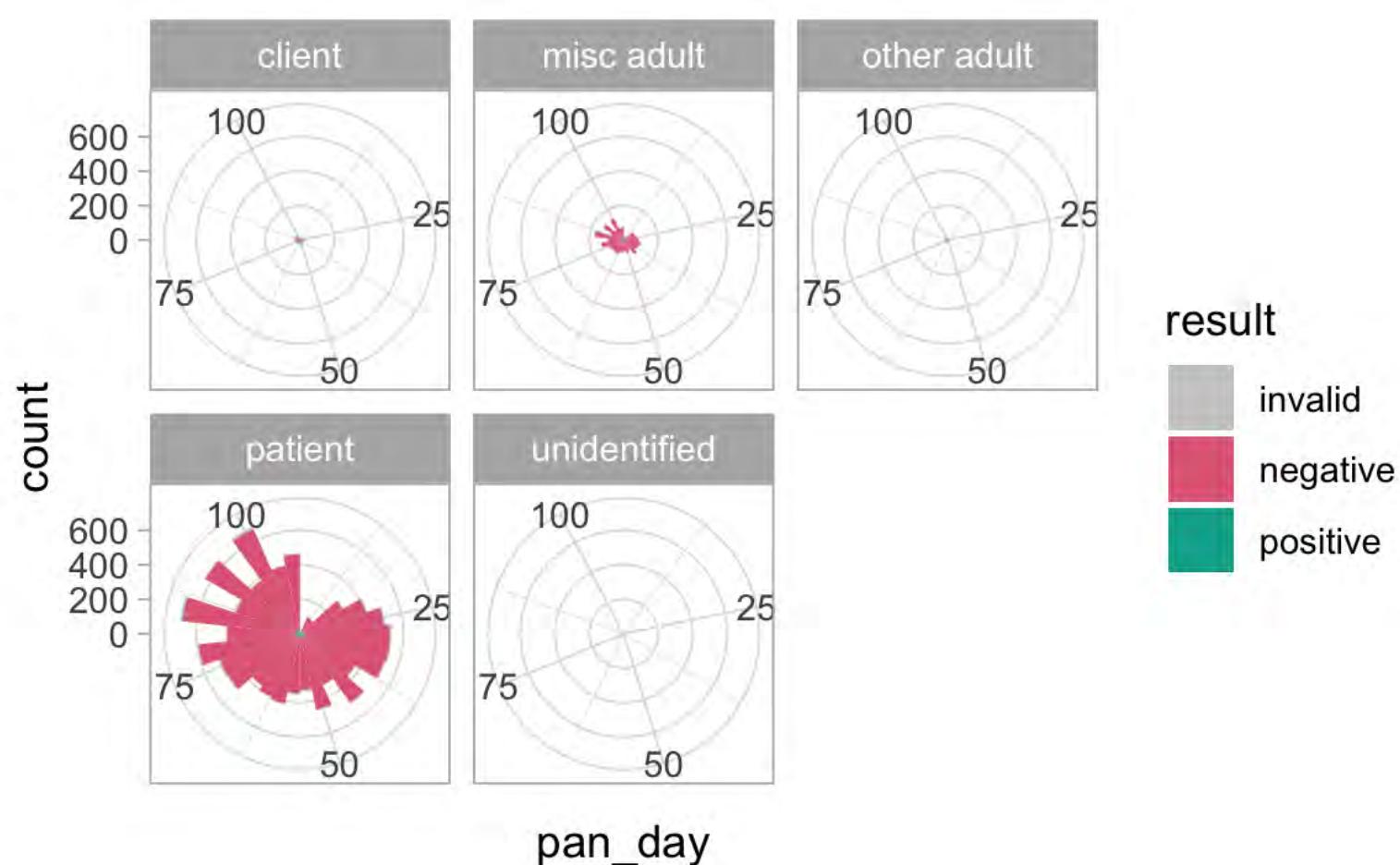


Facets

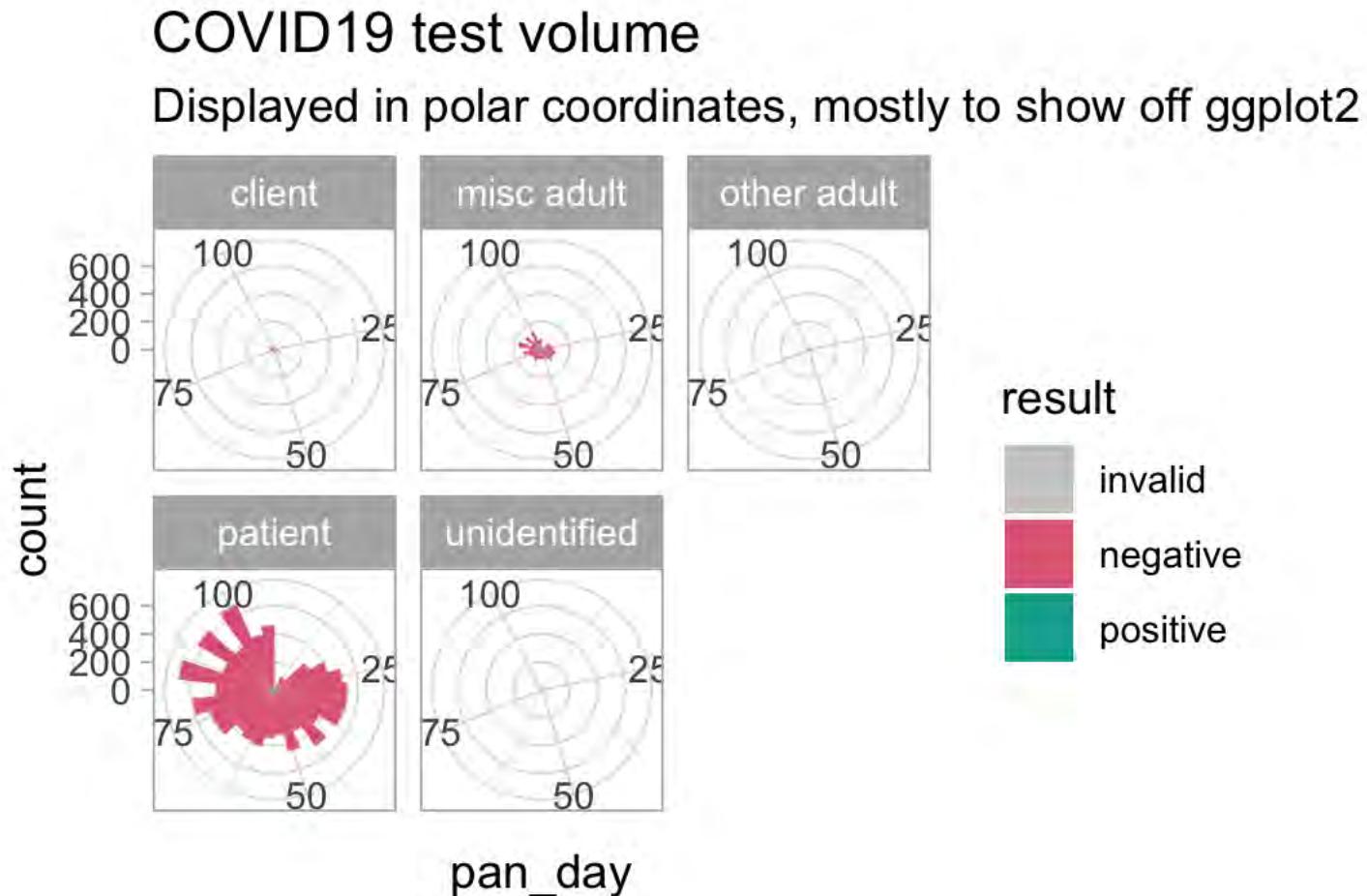
Subplots that display subsets of the data



Coordinate systems



Titles and captions



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



Required

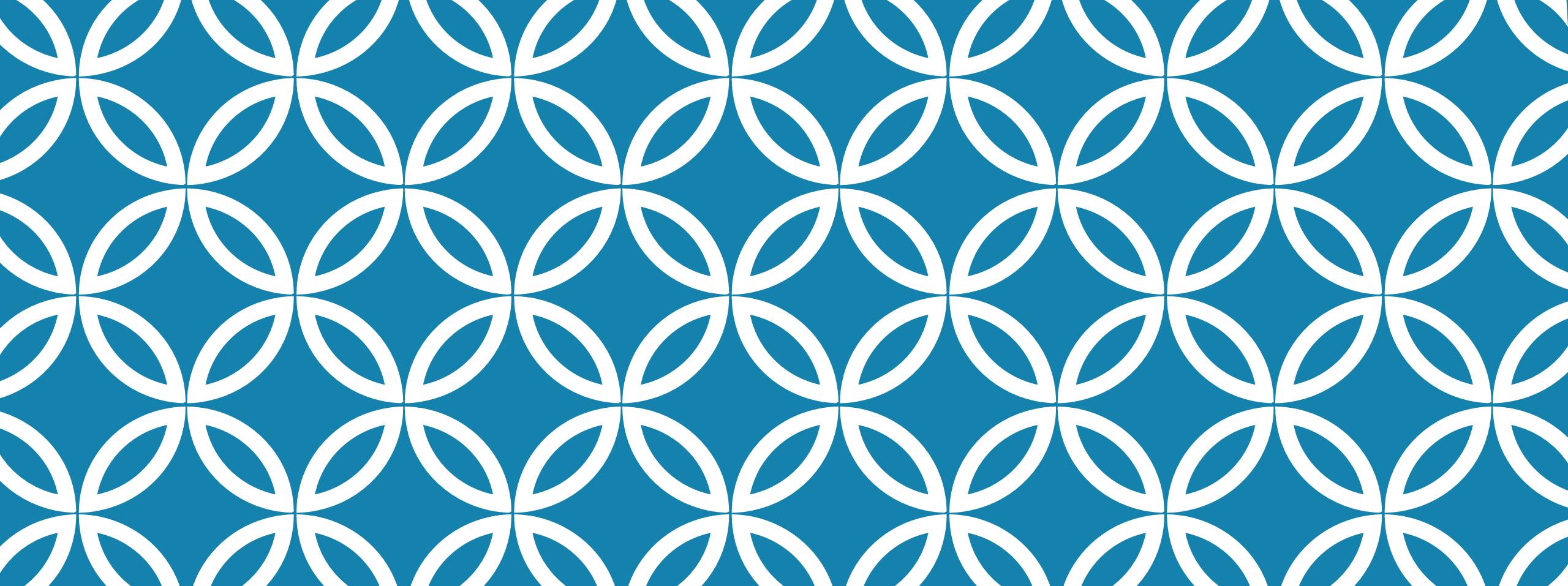
Optional

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: Slicing Your Data

Patrick Mathias
December 13, 2020

Presentation adapted from...

Amrom Obstfeld

Assistant Professor of Clinical Pathology
and Laboratory Medicine

University of Pennsylvania Perelman
School of Medicine

Director of Hematology and Coagulation
Laboratories

Children's Hospital of Philadelphia



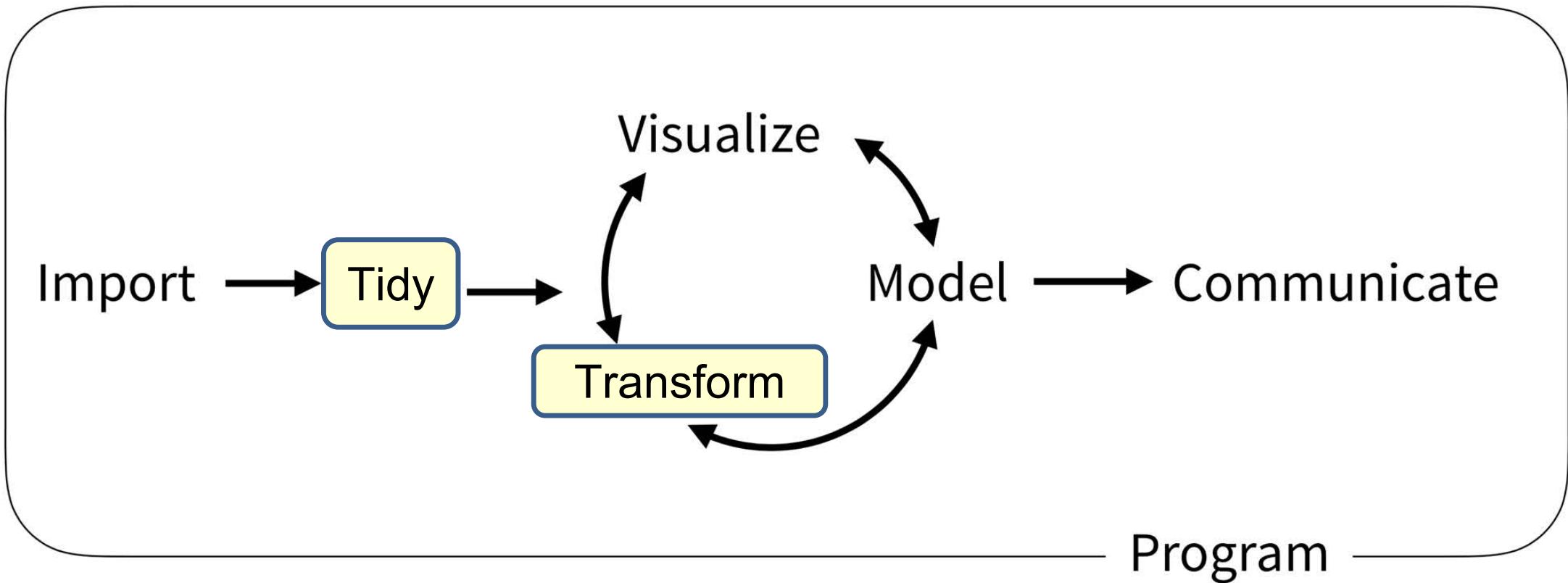
Goal

1. Learn how to use dplyr to transform data frames

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

Typical Data Science Pipeline



What is a “Tidy” Data Frame

A data set is **tidy** if:

AGE	MRN	SEX	RESULT
1	1000000000	M	Normal
2	1000000001	F	Abnormal
3	1000000002	M	Normal
4	1000000003	F	Abnormal

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")
```
```



Transform Data with



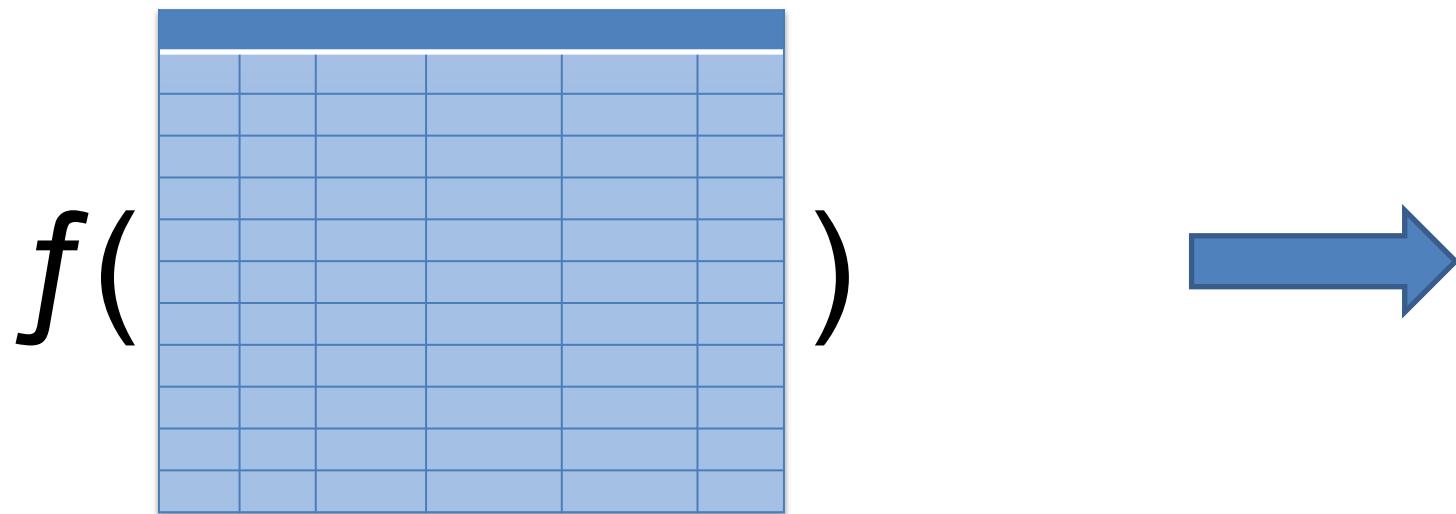
dplyr



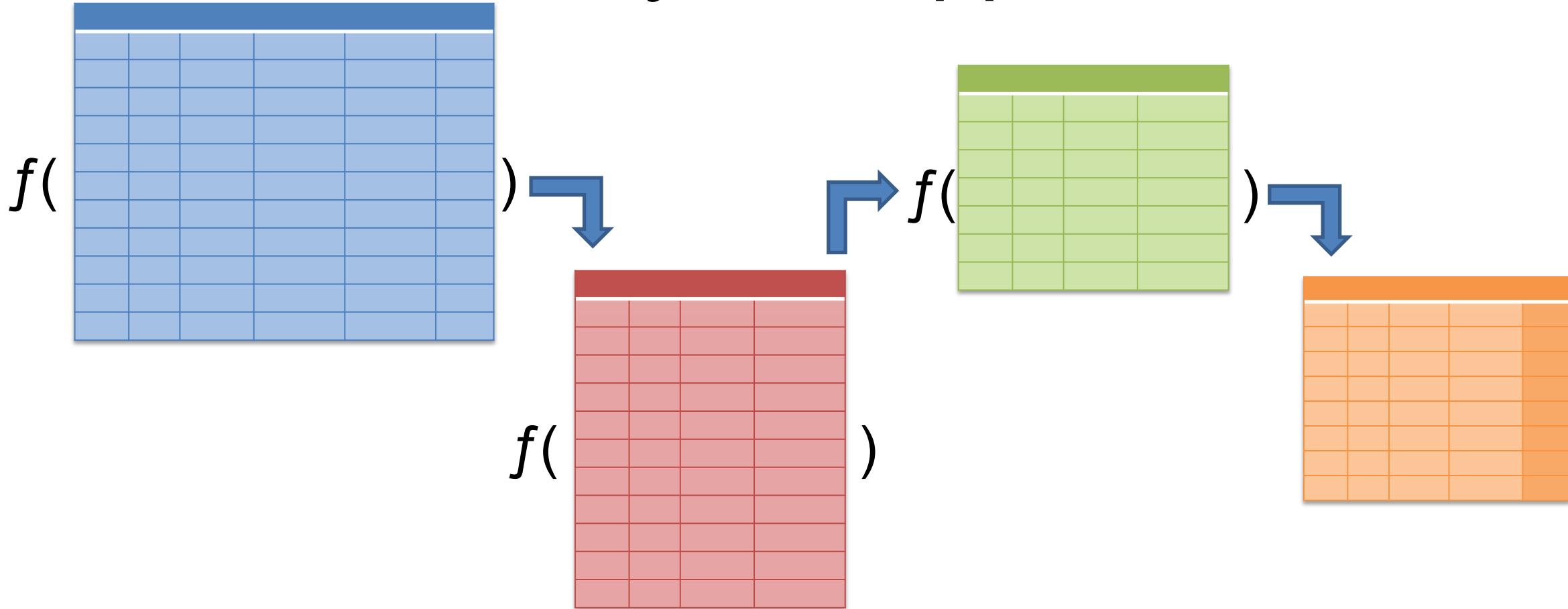
dplyr implements a *grammar* for transforming tabular data.



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.

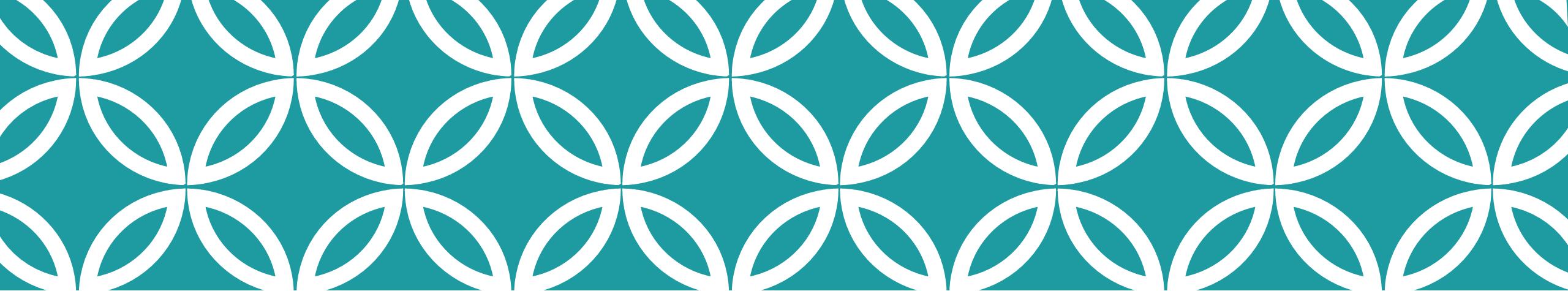
```
function(data, ...)
```

dplyr
function

data frame to
transform

specific
arguments





**Pulling specific columns out of your data
frame**

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 |



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, _____)
```

select() helpers

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)

Manipulate Cases

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

- filter(data, ...)
- distinct(..., keep_all = FALSE)
- sample_frac(tbl, size = 1, replace = FALSE, weight = NULL, ...)
- sample_n(tbl, size, replace = FALSE, weight = NULL, ...)
- slice(...)
- top_n(..., n, wt)

Manipulate Variables

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

- pull(data, var = 1)
- select(data, ...)

VARIATIONS

summarise_all() - Apply funs to every column.
summarise_att() - 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)
- ungroup(x,...)

Logical and boolean operators to use with filter()

- <
- <=
- is.na()
- %in%
- |
- xor()

See ?base::logic and ?Comparison for help.

ARRANGE CASES

- arrange(data, ...)

ADD CASES

- add_row(data, ..., before = NULL, after = NULL)

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)**

R Studio

RStudio® is a trademark of RStudio, Inc. • CC BY SA RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more with browseVignettes(package = c("dplyr", "tidyverse")) • dplyr 0.7.0 • Table 1.20 • Updated: 2017-03





Pulling specific rows out of your data frame

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?

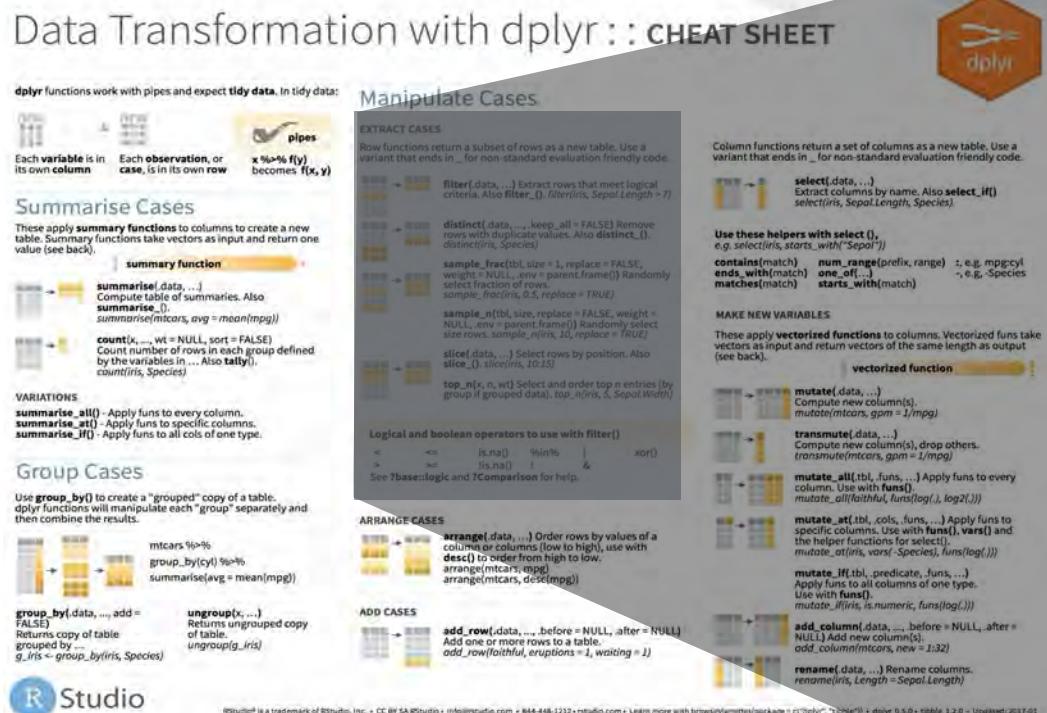
$$3 != 1$$

Your Turn 3

Use filter() with the logical operators to find:

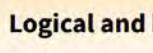
- Every test for patients **over age 80**
- All of the covid testing where the demographic group (demo_group) is **equal to "client"**
- CHALLENGE:
 - All of the covid testing where the patient class (patient_class) **is NA** [Hint: See slide titled "Logical Tests"]

filter() variants



EXTRACT CASES

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

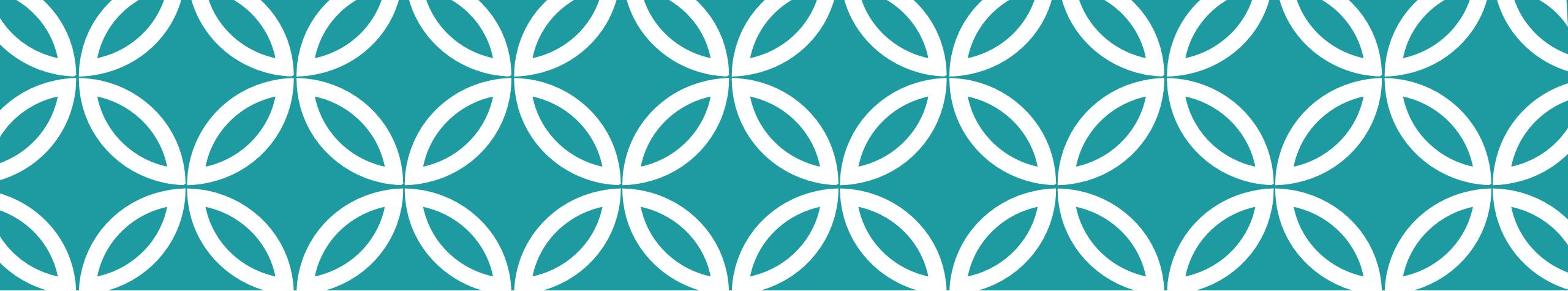
-  **filter(.data, ...)** Extract rows that meet logical criteria. `filter(iris, Sepal.Length > 7)`
-  **distinct(.data, ..., .keep_all = FALSE)** Remove rows with duplicate values. `distinct(iris, Species)`
-  **sample_frac(tbl, size = 1, replace = FALSE, weight = NULL, env = parent.frame())** Randomly select fraction of rows. `sample_frac(iris, 0.5, replace = TRUE)`
-  **slice(.data, ...)** Select rows by position. `slice(iris, 10:15)`
-  **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`. `mutate_all(faithful, funs(log10, log2))`
-  **mutate_at(tbl, cols, funs, ...)** Apply funs to specific columns. Use with `funs`, `cols` and the helper functions for `select`. `mutate_at(iris, c(Sepal.Length, Sepal.Width), funs(log10))`
-  **mutate_if(tbl, predicate, funs, ...)** Apply funs to all columns of one type. Use with `funs`. `mutate_if(iris, is.numeric, funs(log10))`
-  **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()

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

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



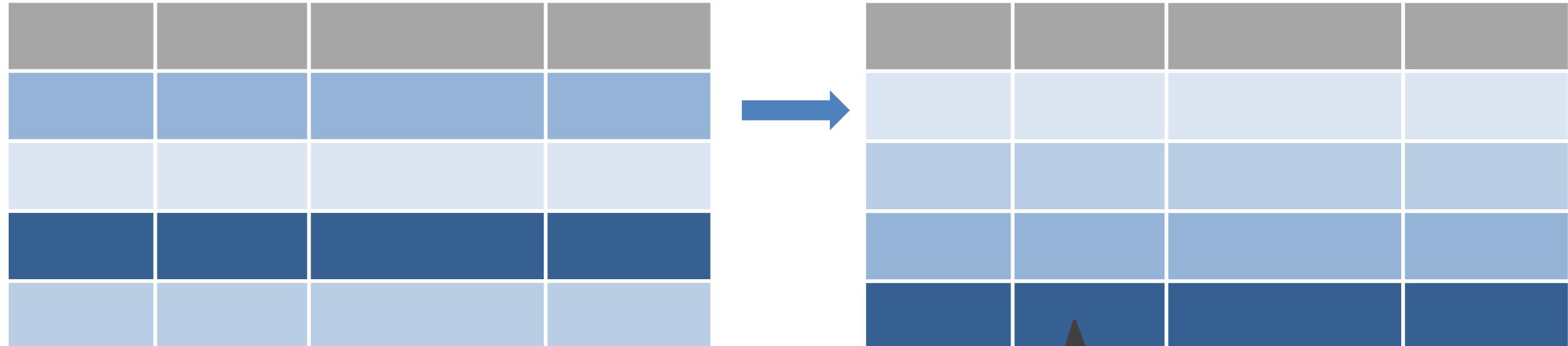


What else?

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_result` 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"?

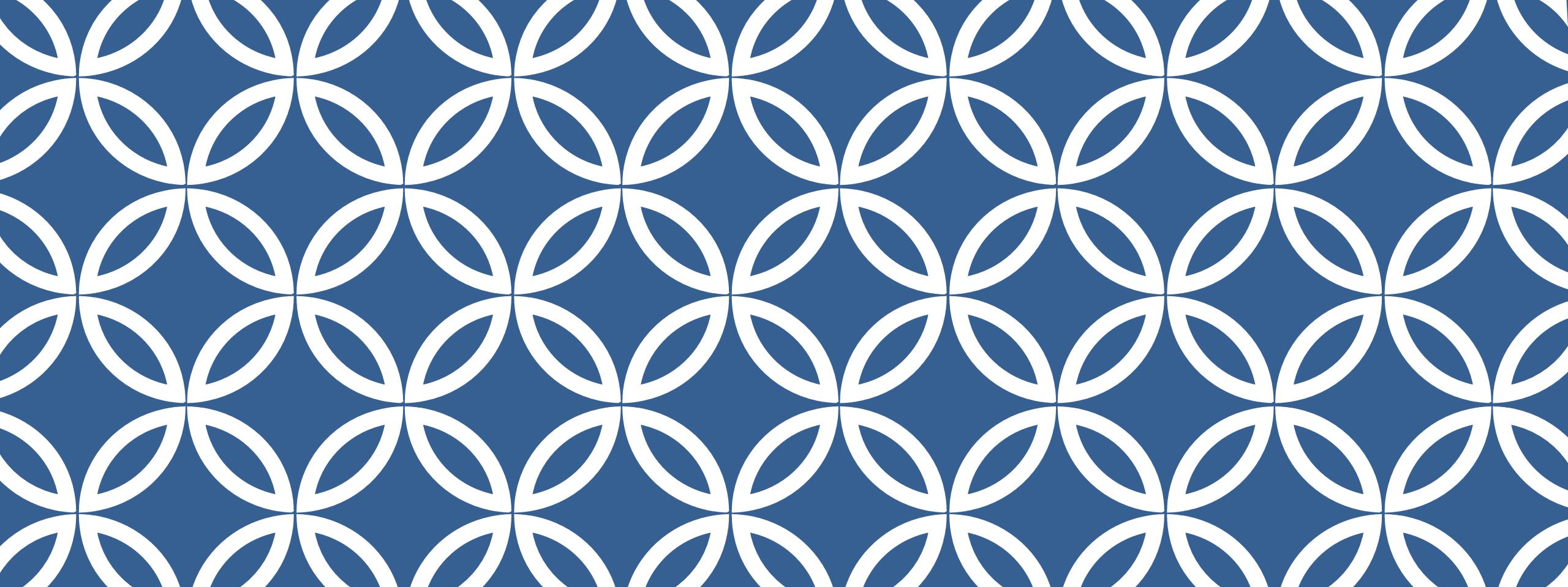


Goal

1. Learn how to use dplyr to transform data frames

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



Data Transformation

Session 5
Dan Herman
December 13, 2020

Presentation adapted from...

Amrom Obstfeld

Assistant Professor of Clinical Pathology
and Laboratory Medicine

University of Pennsylvania Perelman
School of Medicine

Director of Hematology and Coagulation
Laboratories

Children's Hospital of Philadelphia



Goal

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

Objectives

1. Use the pipe operator to pass the output of one function as an input to the next function
2. Create new calculated columns not found in the original data frame

%>%

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 from pandemic day 10 or earlier
2. Select the column that contains the ordering location
3. Arrange those columns by location



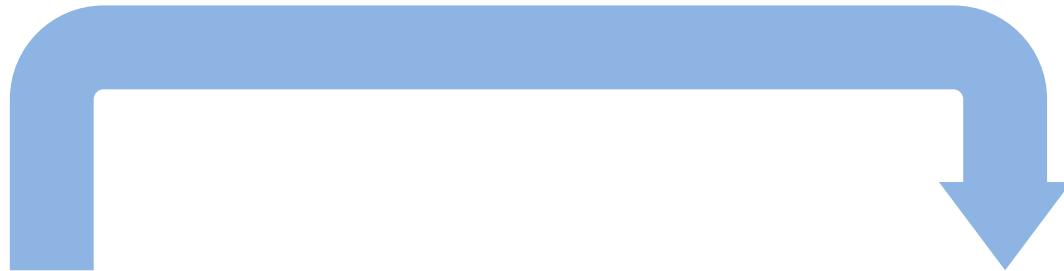
Data Analysis Steps

```
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_testing, pan_day <= 10)
```

```
covid_testing %>% 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)
```



Example: TAT investigation

The PICU would like a word with you because of a recent incident involving a delay in results for a patient who required an aerosol generating procedure.

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



You decide to investigate...WITH DATA

Your Turn 1

Open '05– Transform.Rmd' and run the setup chunk.

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

1. Filters to orders from the clinic (**clinic_name**) of "picu"
2. Selects the columns with the receive to verify turnaround time (**rec_ver_tat**) and day from start of the pandemic (**pan_day**)
3. Arrange the `**pan_day**` from highest to lowest
4. Use <- to assign the result to a new variable and take a look.

Isolating data



Extract variables with `select()`

Extract rows with `filter()`

Arrange rows, with `arrange()`.



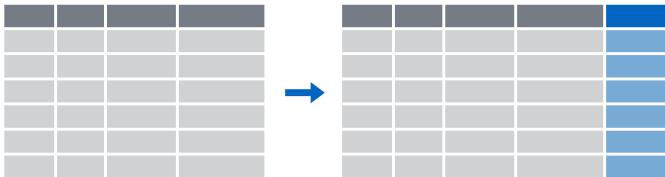
Deriving Data

What are the **average** and
standard deviation of the
total turnaround time for
each ordering clinic?

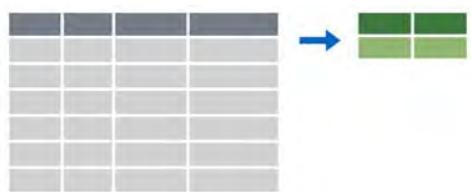
Breaking down the analytical question

1. Total TAT for each order
2. Group orders by clinic
3. Calculate average and standard deviation for each clinic

Deriving data



Make new variables with **mutate()**



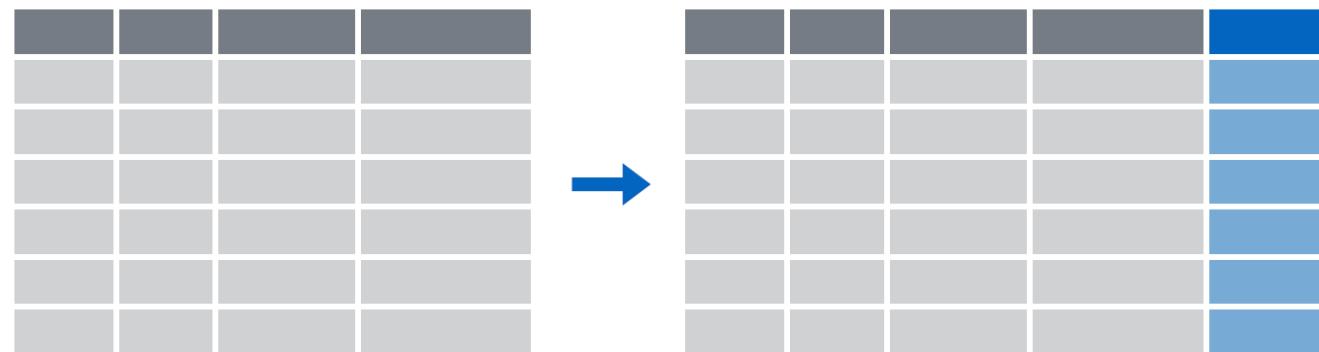
Make summaries of data with
summarize()



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 2

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

mutate()

Replace columns

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

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

| 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()

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 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"

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()
- 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

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 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"

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()
- 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 unique values
- dplyr::sum(is.na()) - # of non-NAs

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
- sdl() - 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 column.
- column_to_rownames() - Move col in row names.
- column_to_rownames(a, var = "C") - Also has_rownames(), remove_rownames()

Combine Tables

COMBINE VARIABLES

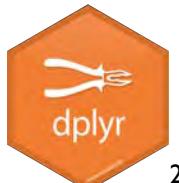
Use **bind_cols()** to paste tables beside each other as a single table.

COMBINE CASES

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

dpqr

on back

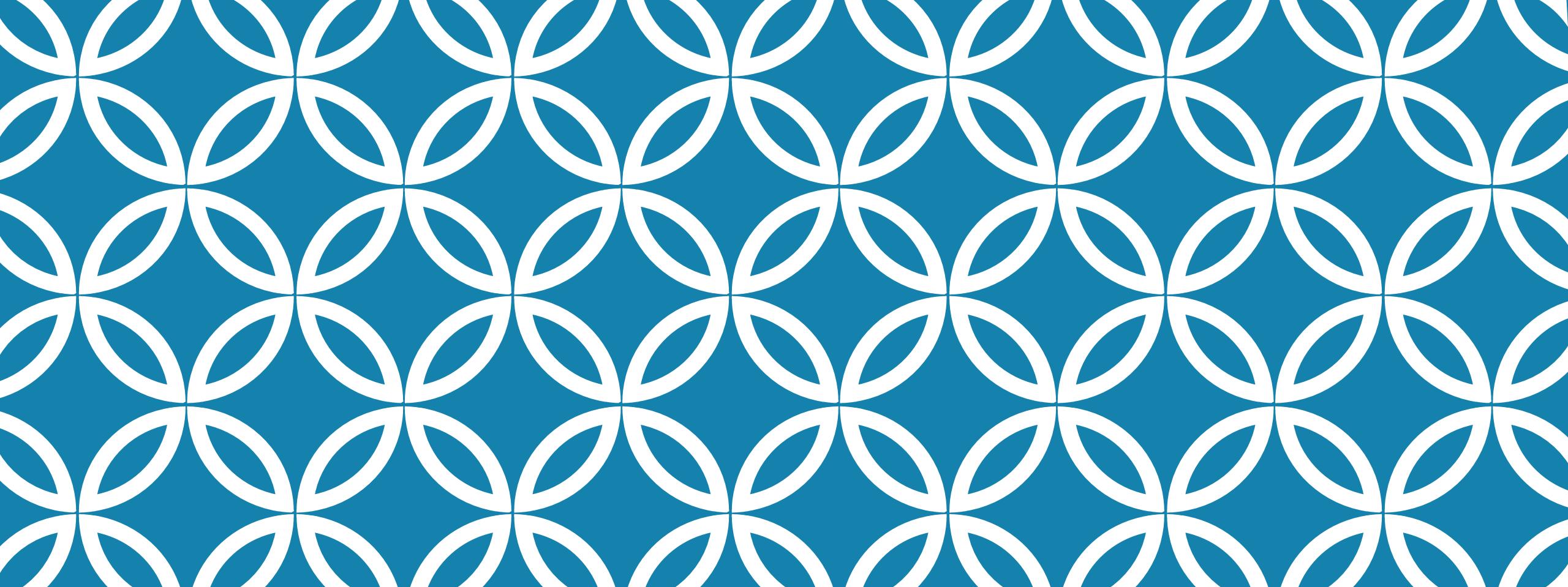


Goal

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

Objectives

1. Use the pipe operator to pass the output of one function as an input to the next function
2. Create new calculated columns not found in the original data frame



Statistical Analyses

Session 6
Dan Herman
December 13, 2020

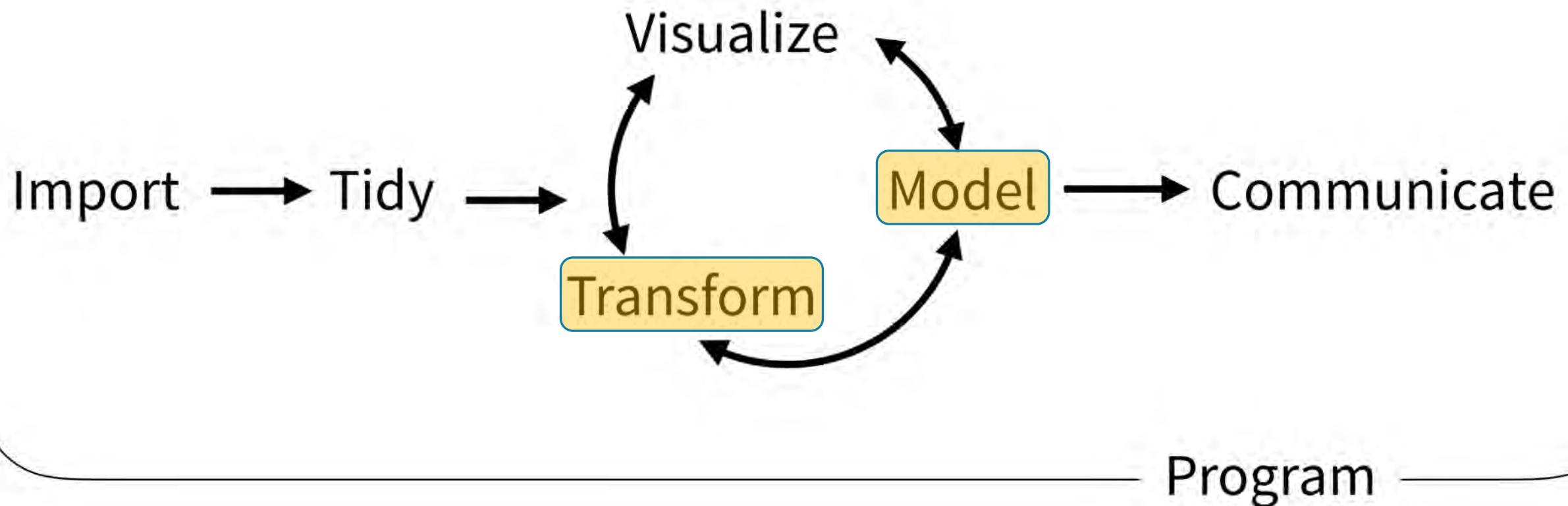
Goals

1. Learn how to summarize data and assess hypotheses

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. Perform a simple test for association

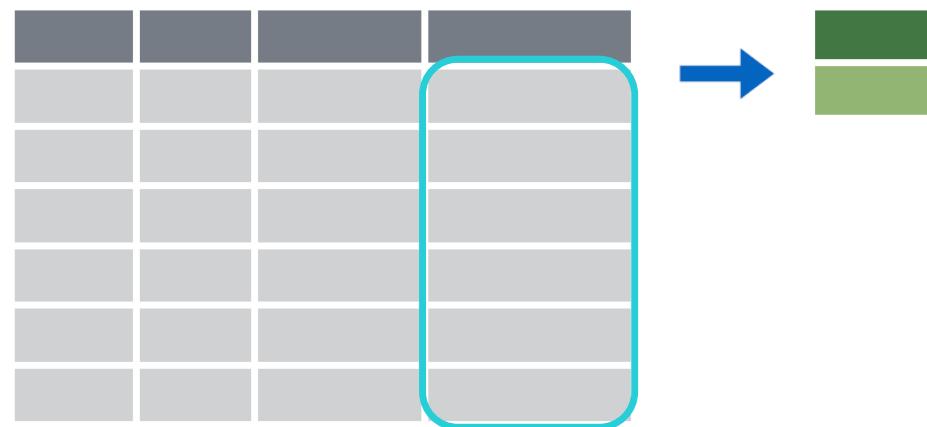
Typical Data Science Pipeline



Summarize()

summarize()

- Make summaries of your data



summarize()

- Make summaries of your data

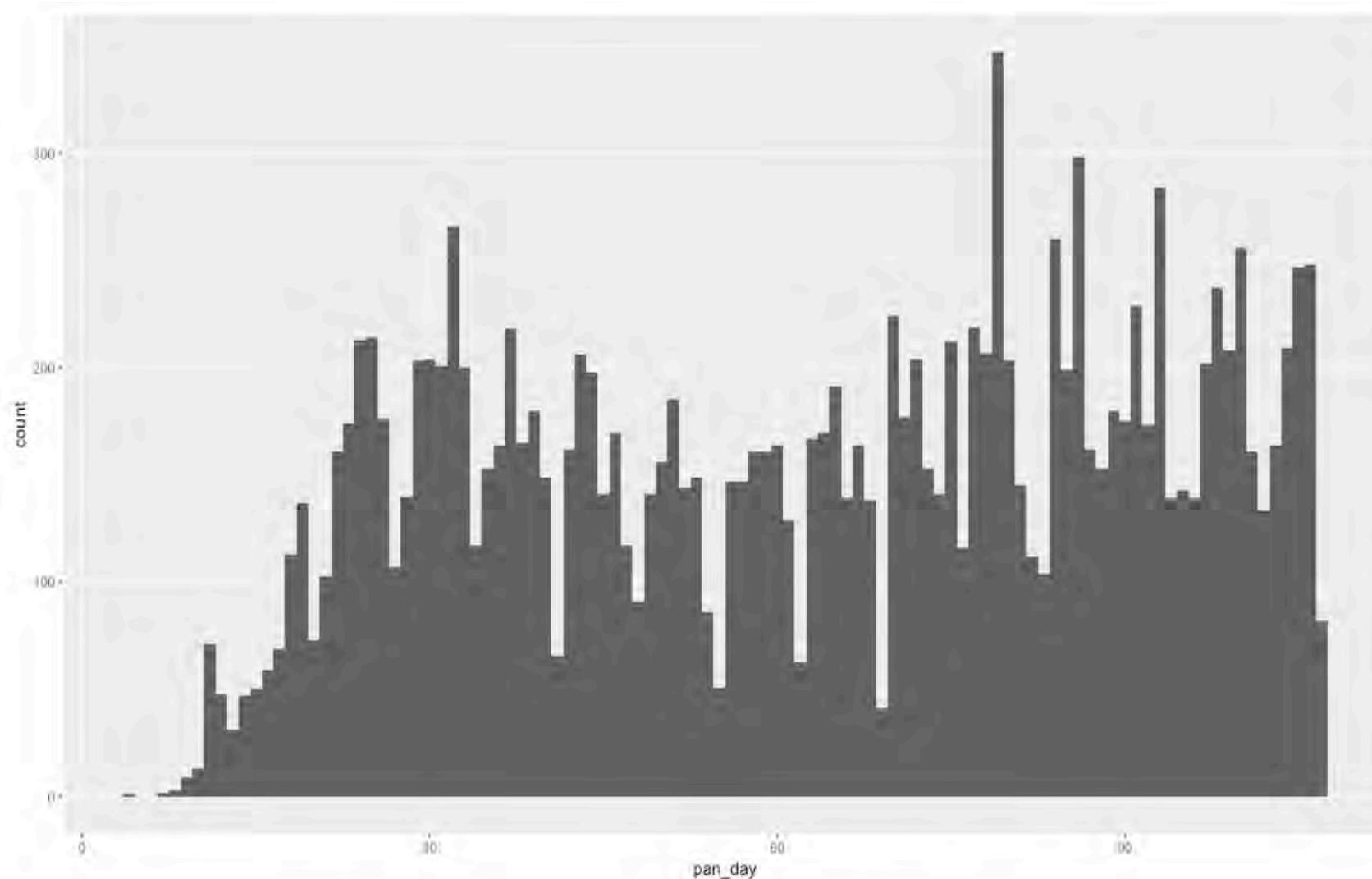
```
covid_testing %>%  
  summarize(new_variable = calculation)
```

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 “06 – 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 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.

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

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

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

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

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

`ABC` 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"))`

`ABC` 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 B C |
| a t 1 | a t 3 |
| b u 2 | b u 2 |
| c v 3 | d w 4 |

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

`ABC` `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)

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

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

`ABC` `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 C |
| a t 1 | a t 3 |
| b u 2 | b u 2 |
| c v 3 | d w 1 |

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



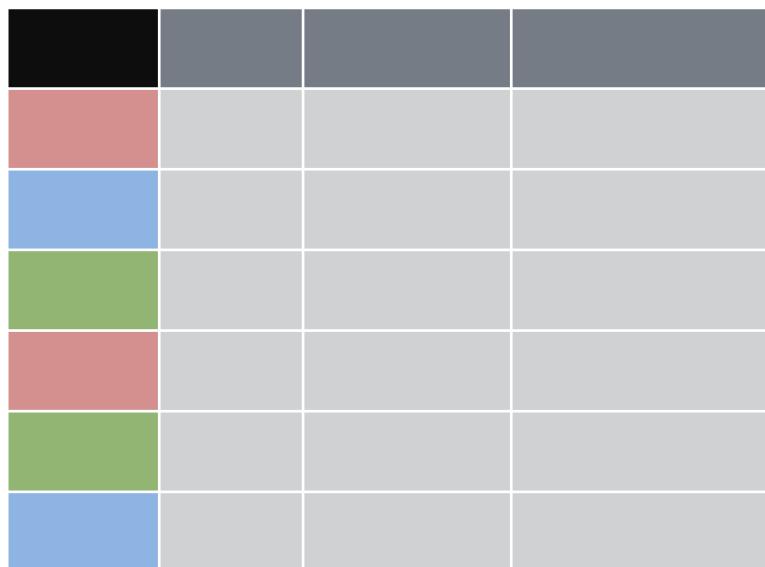
Your Turn 2

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)
```

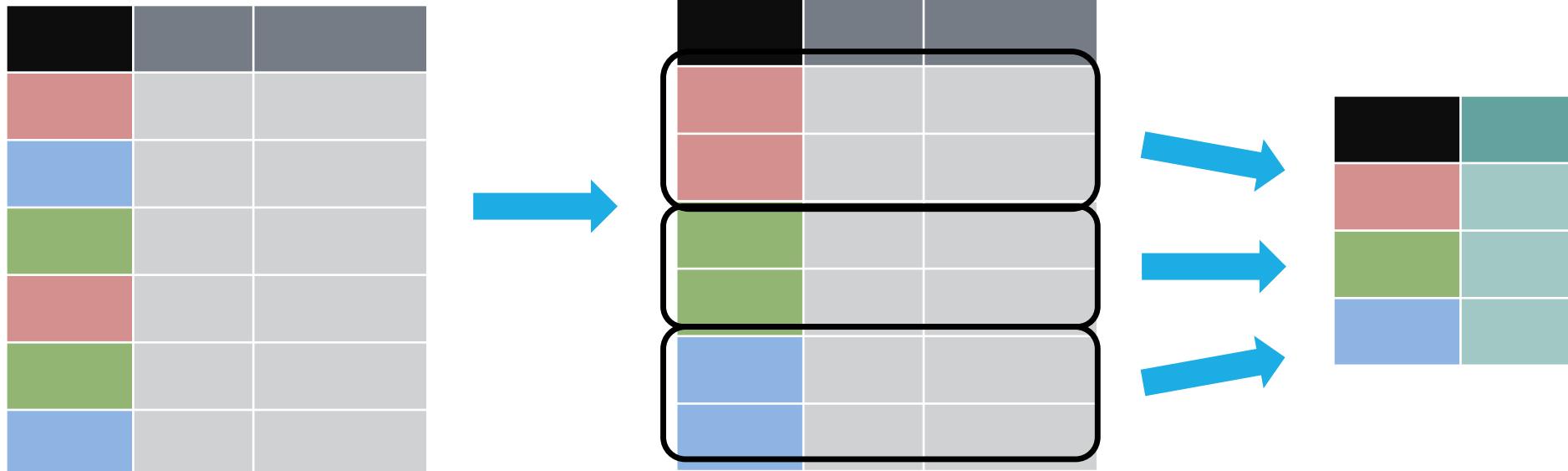
```
# A tibble: 15,524 × 3  
# Groups:   pan_day, clinic_name [2,526]  
  mrn  pan_day clinic_name  
  <dbl>    <dbl>    <chr>  
1 5001412      4 inpatient ward a  
2 5000533      7 clinical lab  
3 5009134      7 clinical lab  
4 5008518      8 clinical lab  
5 5008967      8 emergency dept
```



```
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 |

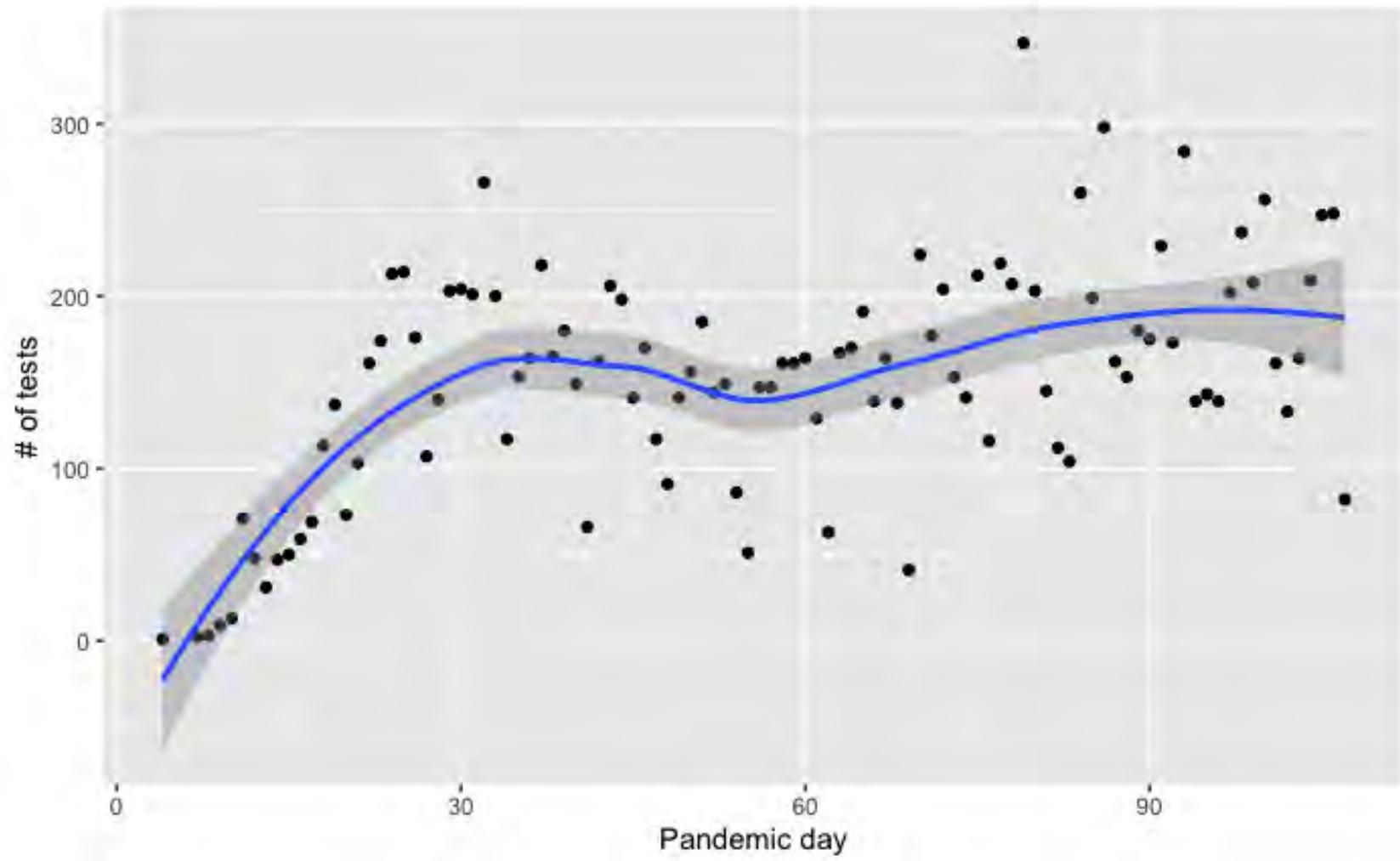


Your Turn 3

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
```



Your Turn 4

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 and assess data

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. Perform a simple test for association

Appendix

R Markdown :: CHEAT SHEET

What is R Markdown?

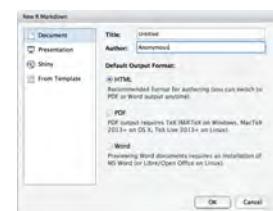


.Rmd files • An R Markdown (.Rmd) file is a record of your research. It contains the code that a scientist needs to reproduce your work along with the narration that a reader needs to understand your work.

Reproducible Research • At the click of a button, or the type of a command, you can rerun the code in an R Markdown file to reproduce your work and export the results as a finished report.

Dynamic Documents • You can choose to export the finished report in a variety of formats, including html, pdf, MS Word, or RTF documents; html or pdf based slides, Notebooks, and more.

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:

- File Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help.
- R Markdown Editor (top):** Shows the YAML header (1), code chunks (2), and various RStudio context menus (3) for code blocks like "set preview location", "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".
- R Markdown Preview (middle):** Displays the rendered HTML content (4). Annotations show "File path to output document" pointing to the header, "synch publish button to accounts at rpubs.com, shinyapps.io" pointing to the "Publish" button, and "RStudio Connect" and "Reload document" buttons.
- R Markdown Console (bottom):** Shows the command `render("report.Rmd", output_file = "report.html")` and its output (5).
- File Explorer (bottom right):** Shows the file structure with `report.Rmd` and `report.html` (6).

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()
```
```

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)

eval - Run code in chunk (default = TRUE)

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)

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

GLOBAL OPTIONS

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

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

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)

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

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**



.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)
})..
```

speed	dist
1	4.00
2	10.00
3	14.00
4	19.00
5	22.00
6	26.00
7	30.00
8	33.00
9	36.00
10	39.00

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/superscript<sup>2</sup>  
~~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}  
<em>HTML ignored in pdfs</em>

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

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

$E = mc^2$

block quote

## Header1

## Header 2

### Header 3

#### Header 4

##### Header 5

##### Header 6

HTML ignored in pdfs

<http://www.rstudio.com>

link

Jump to Header 1

image:



Caption

- 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

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<sup>[1]</sup>

1. Here is the footnote.<sup>[2]</sup>

# 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):

Indent 2 spaces      Indent 4 spaces

```

```

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

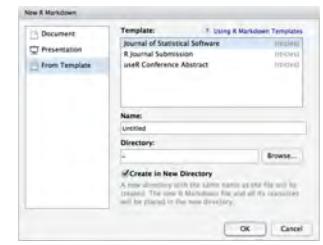
```

```

name: My Template

```

```



### sub-option

### description

		html	pdf	word	odt	rtf	md	gitlab	ioslides	slidy	beamer
citation_package	The LaTeX package to process citations, natbib, biblatex or none	X									
code_folding	Let readers to toggle the display of R code, "none", "hide", or "show"		X								
colortheme	Beamer color theme to use										X
css	CSS file to use to style document		X						X	X	
dev	Graphics device to use for figure output (e.g. "png")		X	X					X	X	X
duration	Add a countdown timer (in minutes) to footer of slides										X
fig_caption	Should figures be rendered with captions?	X	X	X	X				X	X	X
fig_height, fig_width	Default figure height and width (in inches) for document	X	X	X	X	X	X	X	X	X	X
highlight	Syntax highlighting: "tango", "pygments", "kate", "zenburn", "textmate"	X	X	X					X	X	
includes	File of content to place in document (in_header, before_body, after_body)	X	X	X	X	X	X	X	X	X	X
incremental	Should bullets appear one at a time (on presenter mouse clicks)?								X	X	X
keep_md	Save a copy of .md file that contains knitr output		X	X	X	X			X	X	
keep_tex	Save a copy of .tex file that contains knitr output										X
latex_engine	Engine to render latex, "pdflatex", "xelatex", or "lualatex"										X
lib_dir	Directory of dependency files to use (Bootstrap, MathJax, etc.)		X						X	X	
mathjax	Set to local or a URL to use a local/URL version of MathJax to render equations		X						X	X	
md_extensions	Markdown extensions to add to default definition or R Markdown	X	X	X	X	X	X	X	X	X	X
number_sections	Add section numbering to headers		X	X							
pandoc_args	Additional arguments to pass to Pandoc	X	X	X	X	X	X	X	X	X	X
preserve_yaml	Preserve YAML front matter in final document?										X
reference_docx	docx file whose styles should be copied when producing docx output										X
self_contained	Embed dependencies into the doc										X
slide_level	The lowest heading level that defines individual slides										X
smaller	Use the smaller font size in the presentation?										X
smart	Convert straight quotes to curly, dashes to em-dashes, ... to ellipses, etc.	X								X	X
template	Pandoc template to use when rendering file quarterly_report.html	X	X	X						X	X
theme	Bootswatch or Beamer theme to use for page	X									X
toc	Add a table of contents at start of document	X	X	X	X	X	X	X	X	X	X
toc_depth	The lowest level of headings to add to table of contents	X	X	X	X	X	X	X	X	X	X
toc_float	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	
<code>eruptions waiting</code>	
1	3.60 79.00
2	1.80 54.00
3	3.33 74.00
4	2.28 62.00
2.283	62

eruptionswaiting	
1	3.600 79
2	1.800 54
3	3.333 74
4	2.283 62

Table with xtable

# 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

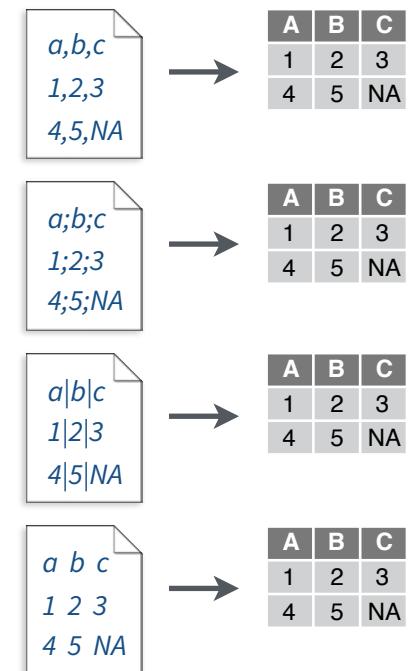
```
write_tsv(x, path, na = "NA", append = FALSE,
 col_names = !append)
```



## 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))
```

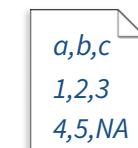
```
write_file(x = "a b c\n1 2 3\n4 5 NA", path = "file.fwf")
```

### 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

<b>tibble(...)</b>	Construct by columns. <code>tibble(x = 1:3, y = c("a", "b", "c"))</code>	Both make this tibble
<b>tribble(...)</b>	Construct by rows. <code>tribble(~x, ~y, 1, "a", 2, "b", 3, "c")</code>	<code>A tibble: 3 × 2</code> <code>x &lt;int&gt; y &lt;chr&gt;</code> 1 1 a 2 2 b 3 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 **tidyr**

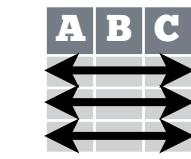
**Tidy data** is a way to organize tabular data. It provides a consistent data structure across packages.

A table is tidy if:



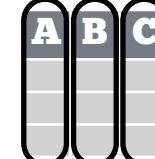
Each **variable** is in its own **column**

&



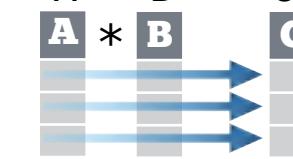
Each **observation**, or **case**, is in its own **row**

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")`

**spread(data, key, value, fill = NA, convert = FALSE, drop = TRUE, sep = NULL)**

**spread()** moves the unique values of a **key** column into the column names, spreading the values of a **value** column across the new columns.

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	2000	cases	80K
B	1999	pop	172M
B	2000	cases	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
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

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

country	year
Afghan	1999
Afghan	2000
Brazil	1999
Brazil	2000
China	1999
China	2000

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

## Handle Missing Values

**drop\_na(data, ...)**

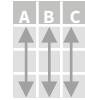
Drop rows containing NA's in ... columns.

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

# 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.

**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



## 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.