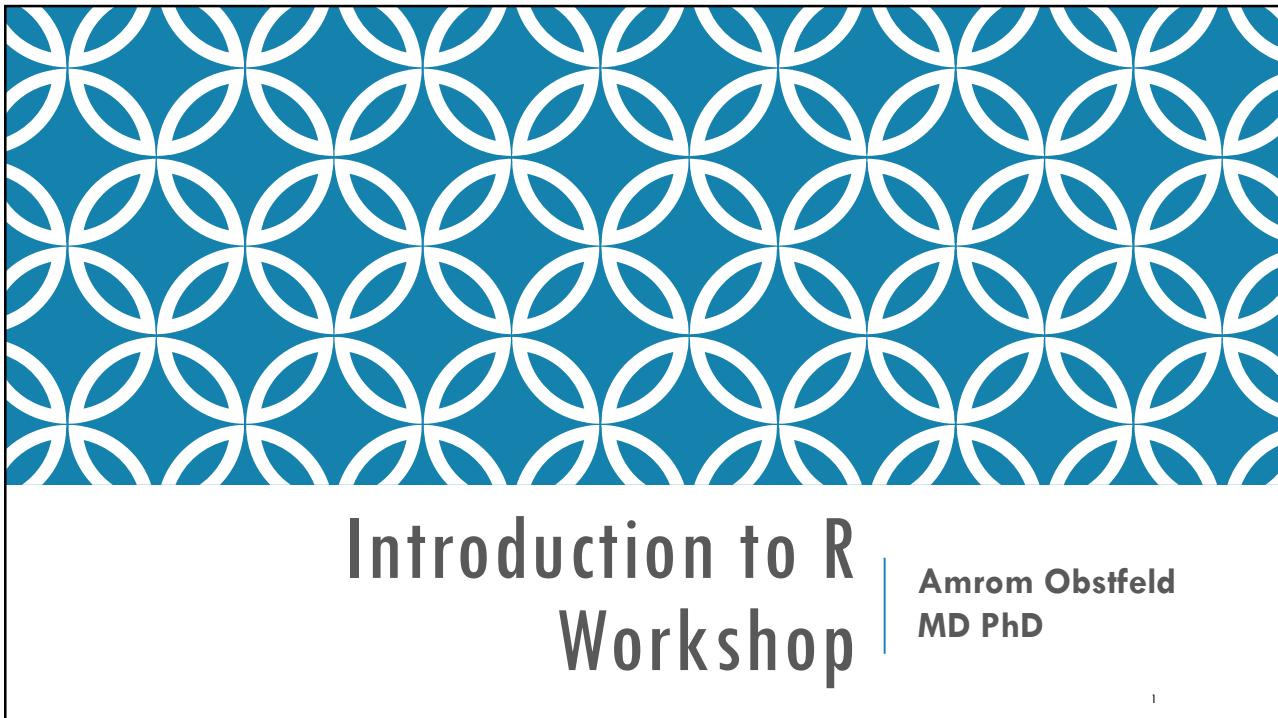


# API Introduction to R Workshop

May 9, 2022



1



2

## Goals and Objectives

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

3

3

| Session  | Instructor      |
|--|-----------------|
| Instructor Introductions, Introduction to technology | Amrom Obstfeld  |
| Introduction to R and RStudio                        | Joe Rudolf      |
| Reproducible Reporting                               | Joe Rudolf      |
| Data Visualization                                   | Stephan Kadauke |
| <hr/>  |                 |
| Data Transformation                                  | Amrom Obstfeld  |
| Group and Summarize                                  | Patrick Mathias |
| Advanced Reporting                                   | Patrick Mathias |

4

4

# Who are we?

5

## Joseph Rudolf

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

Medical Director, Automated Core Laboratory, ARUP  
Laboratories



6

6

## Patrick Mathias

Assistant Professor, Department of  
Laboratory Medicine and Pathology

University of Washington School of  
Medicine

Associate Medical Director, Laboratory  
Medicine and Pathology Informatics



7

## 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



8

8

## Amrom Obstfeld

Assistant Professor of Clinical Pathology  
and Laboratory Medicine

University of Pennsylvania Perelman  
School of Medicine

Director of Pathology Informatics

Children's Hospital of Philadelphia

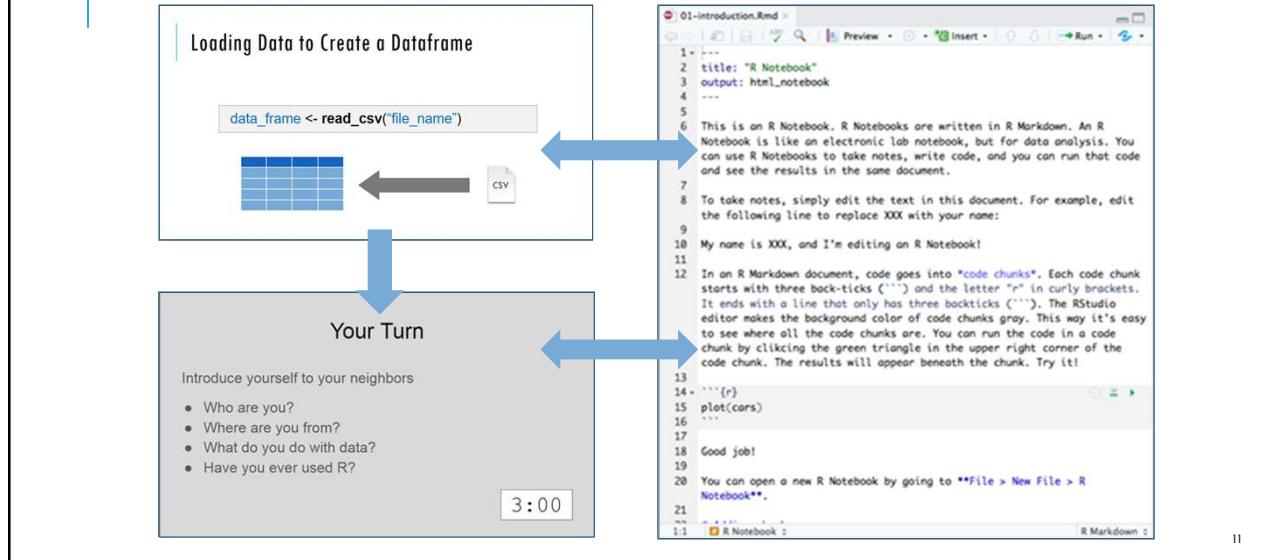


9

## Workshop Workflow

10

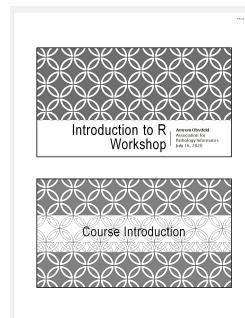
# Sessions



11

# Workshop Coursebook

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



12

12



## Who are you?

13

## Your Turn

Introduce yourself to your breakout roommates

Who are you?  
Where are you from?  
Why are you here?  
Have you ever used R?

14

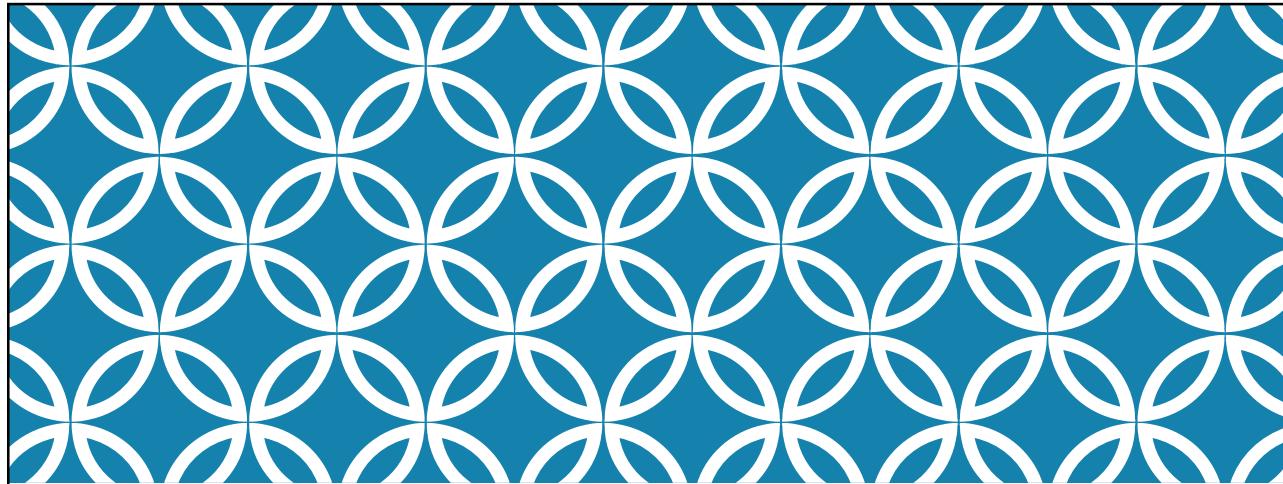
14

## Final Tips

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

15

15



# **Introduction to R, RStudio, and R Markdown**

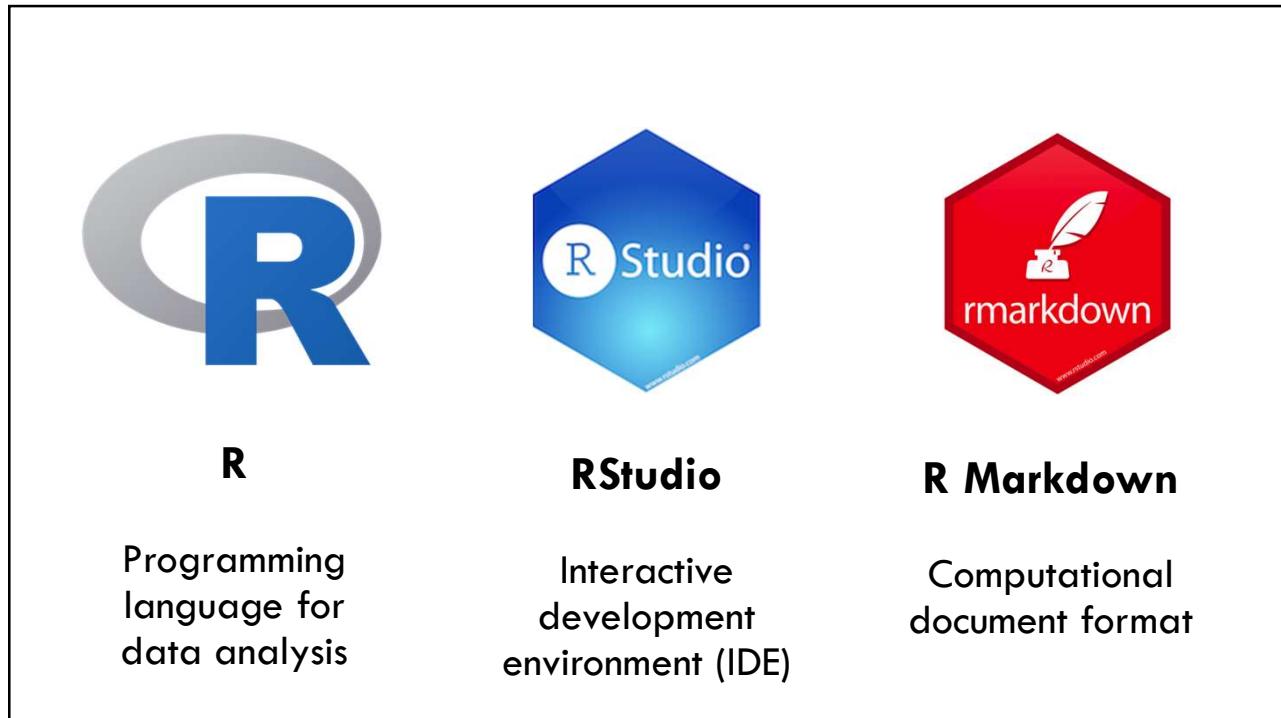
Joseph Rudolf  
API-R 2022

1

## **Part I**



2



3



4

## RStudio: On the Web and In Your Home



**RStudio Server**  
Hosted on a server  
(in the cloud)



**RStudio Desktop**  
Installed locally on  
your computer

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

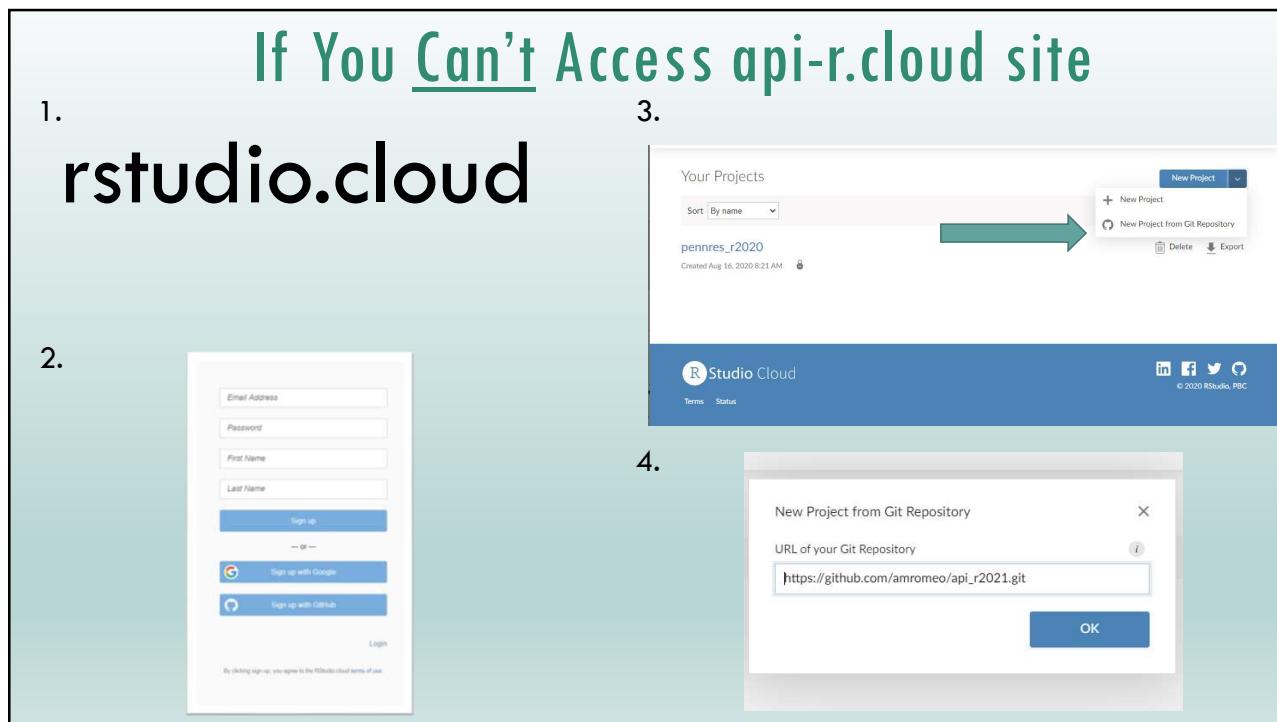
5

## Your Turn #1

Go to <https://api-r.cloud> in your browser and log in using the username and password provided in the course email.

Click “thumbs up” in zoom once you see the RStudio panes.

If you can't access the site click “thumbs down” and we will set you up in a backup configuration shortly.



7

**EDITOR**

```
1+ ---
2 title: ''
3 output:
4
5+ ---
6
7+ # Heading
8 |
9
10+ ``{r}
11+
12 Code
13
14+ ```
15
16
8:1  Heading 2
```

**CONSOLE**

```
R version 4.0.2 (2020-06-22) -- "Taking Off Again"
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.
```

**ENVIRONMENT**

Environment is empty

**MISC**

8

# Reproducible Data Analysis and R Markdown

9

## The Duke Cancer Scandal

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



10

Duke

MD Anderson

|              |              |
|--------------|--------------|
| "1881_at"    | "1882_g_at"  |
| "31321_at"   | "31322_at"   |
| "31725_s_at" | "31726_at"   |
| "32307_r_at" | "32308_r_at" |

...

**Off-by-one indexing error**

11

“Common problems are simple...

Off-by-one **indexing error**

Sensitive / resistant **label reversal**

**Confounding** in experimental design

Inclusion of data from **non-reported sources**

**Wrong figure** shown

... and simple problems are common.”

12

## Point-and-click is not reproducible



Computer code can precisely document each step of the analysis

13

## Why YOU should analyze your data 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 from my analysis?”



**YOUR CLOSEST COLLABORATOR IS **YOU** FROM 6 MONTHS AGO**

14



15

## Anatomy of an R Markdown Document

```
1 ---  
2 title: 'My Markdown Document'  
3 output: html_document  
4 ---  
5  
6 # One Hashtag = Large Header  
7  
8 ## Two Hashtags = Smaller Header  
9  
10 Here is some text.  
11  
12 * It's easy to make a list  
13 * Here is how you style text cursive or bold  
14  
15  
16 ``{r}  
17 x <- rnorm(100)  
18 summary(x)  
19 ``
```

Header

Text  
(with marks)

Code chunk

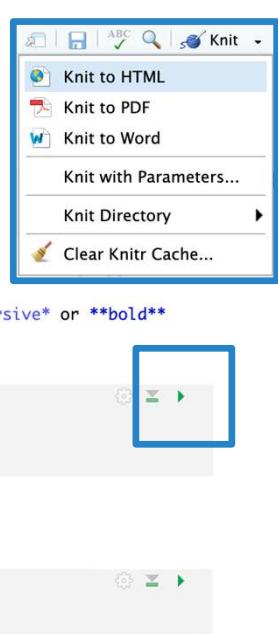
16

1 `---`  
 2 `title: 'My Markdown Document'`  
 3 `output: html_document`  
 4 `---`

5  
 6 `# One Hashtag = Large Header`  
 7  
 8 `## Two Hashtags = Smaller Header`  
 9  
 10 Here is some text.  
 11  
 12 \* It's easy to make a list  
 13 \* Here is how you style text *cursive* or **bold**  
 14

15 ````{r}`  
 16 `x <- rnorm(100)`  
 17 `summary(x)`  
 18 `````

20  
 21 `## Including Plots`  
 22  
 23 ````{r, echo=FALSE}`  
 24 `hist(x)`  
 25 `````



**My Markdown Document**

## One Hashtag = Large Header

## Two Hashtags = Smaller Header

Here is some text.

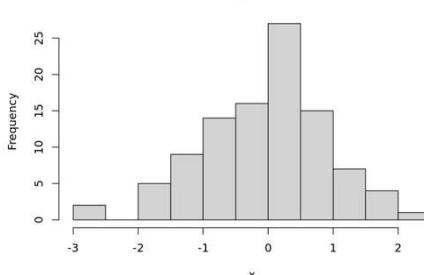
- It's easy to make a list
- Here is how you style text *cursive* or **bold**

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

```
#>    Min. 1st Qu. Median Mean 3rd Qu. Max.
#> -2.99204 -0.64726 0.14853 -0.02832 0.58218 2.07410
```

**Including Plots**

**Histogram of x**



17

1 `---`  
 2 `title: 'My Markdown Document'`  
 3 `output: html_document`  
 4 `---`

5  
 6 `# One Hashtag = Large Header`  
 7  
 8 `## Two Hashtags = Smaller Header`  
 9  
 10 Here is some text.  
 11  
 12 \* It's easy to make a list  
 13 \* Here is how you style text *cursive* or **bold**  
 14

15 ````{r}`  
 16 `x <- rnorm(100)`  
 17 `summary(x)`  
 18 `````

20  
 21 `## Including Plots`  
 22  
 23 ````{r, echo=FALSE}`  
 24 `hist(x)`  
 25 `````



**My Markdown Document**

## One Hashtag = Large Header

## Two Hashtags = Smaller Header

Here is some text.

- It's easy to make a list
- Here is how you style text *cursive* or **bold**

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

```
#>    Min. 1st Qu. Median Mean 3rd Qu. Max.
#> -2.99204 -0.64726 0.14853 -0.02832 0.58218 2.07410
```

**Including Plots**

**Histogram of x**



18

## Your Turn #2

Open a sample R Markdown document (File -> New File -> R Markdown).

Review the format of the document: header, text, code chunks

Execute the individual code chunks by selecting the Run Current Chunk arrow.

Knit the document to HTML (Preview or Knit Button -> Knit to HTML). You may be prompted to save your R Markdown first. In this case select a name for your document and click save. Review the knitted document.

03:00

19

## Part II

20

# The Basics of Coding

21

## The Basics of Coding: Calculation

- R is a calculator!

```
1
2 ^ ``{r}
3
4
5
6 ^ ``
```

[1] 7

press play  
button to execute  
code

answer returned  
here

22

## The Basics of Coding: Functions

- Code that extends our reach beyond the basic operators

```

1
2 ``{r}
3
4 abs(-77)
5
6 ````

[1] 77
  
```

function (does stuff)

argument (input)

**abs(-77)**

23

## 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

```

1
2 ``{r}
3
4 seq(from=5, to=150, by=10)
5
6 ````

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

24

## The Basics of Coding: Objects

- Objects are the container for your output

**object**  
(stores output)

**function**  
(does stuff)

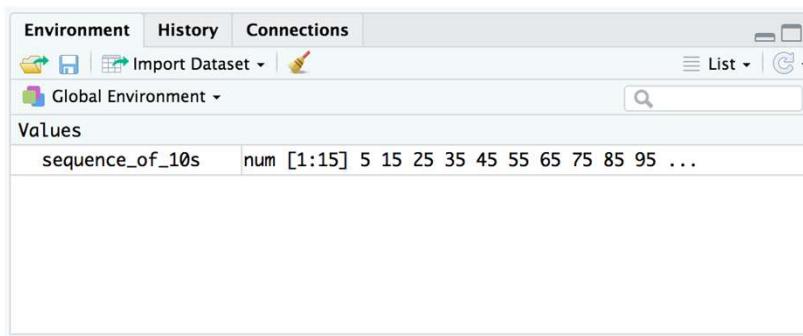
**arguments**  
(input)

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

25

## Checking the contents of an object

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



26

## 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!

The image shows two side-by-side RStudio code editors. Both have identical code:

```

1
2 ~`{r}
3
4 min(sequence_of_10s)
5
6 ~```

```

The left editor shows the output [1] 5. The right editor shows the output [1] 145.

27

## Your Turn #3

I've written some code to create a sequence from 0 to 500 in increments of 25 called `sequence_of_25s`. Ultimately I want to calculate the median value of this sequence. Unfortunately I've made some mistakes in my code and I am hoping you can help me find them.

The image shows an RStudio code editor with the following code:

```

1
2 ~`{r}
3
4 sequence_of_25s <- seq(from=0 to=50, by=25)
5
6 ~```
7
8 ~`{r}
9
10 median(sequence_of_25s)
11
12 ~```
13

```

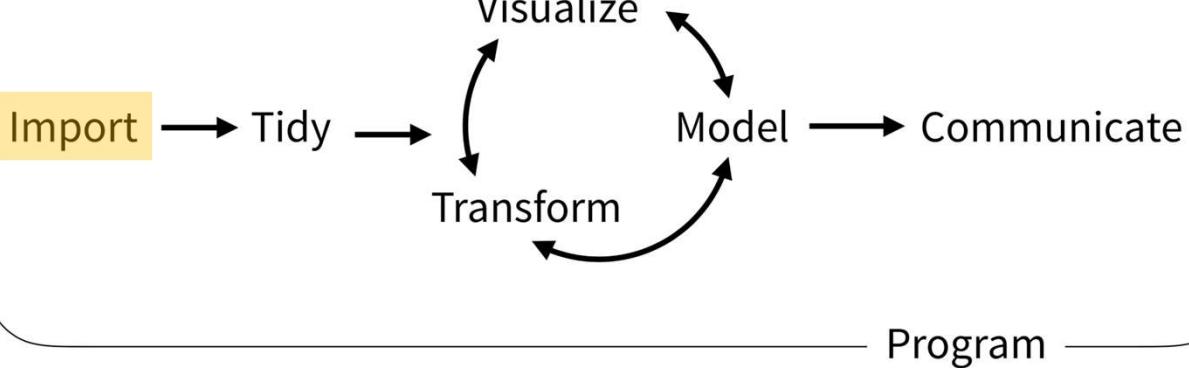
To the right of the code editor is a digital timer displaying 01:00.

28

## Importing Data

29

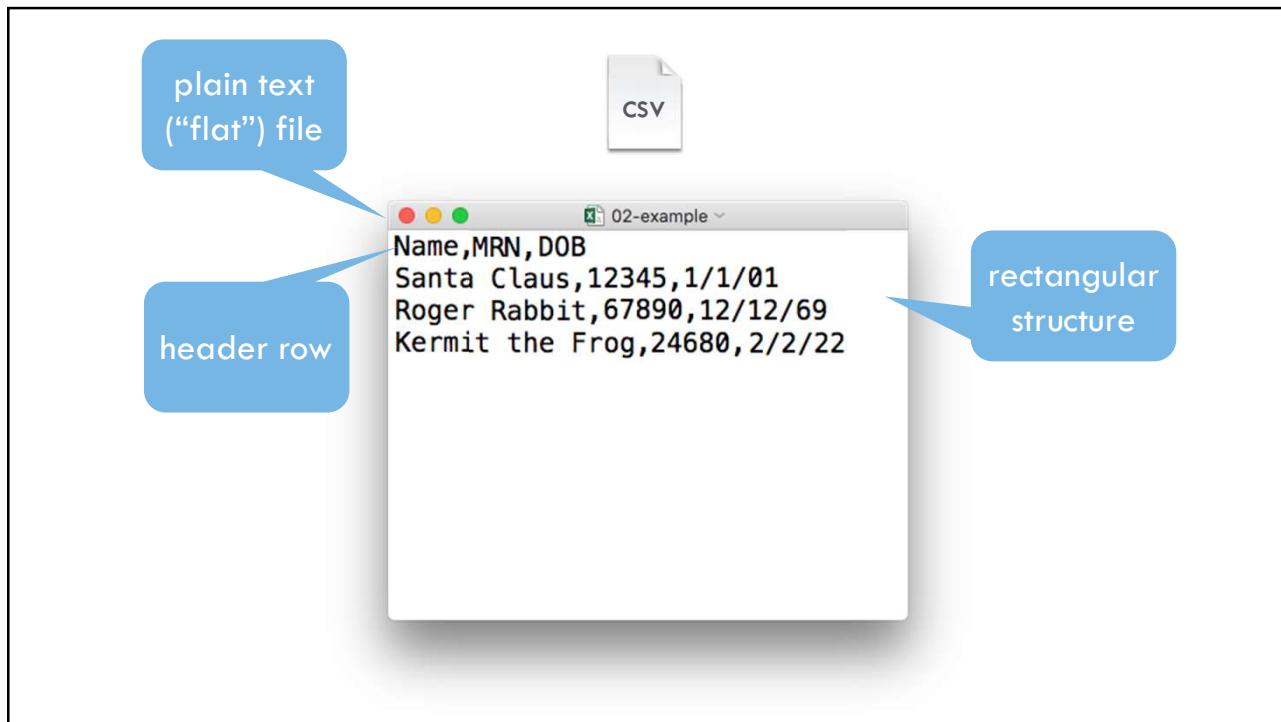
## The Data Analysis Pipeline



From *R for Data Science* (<https://r4ds.had.co.nz/introduction.html>)

30

30



31

## Tidyverse: R Packages for Data Science

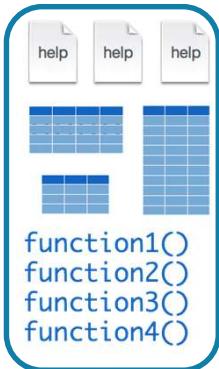
- A consistent way to organize data
- Human readable, concise, consistent code
- Build pipelines from atomic data analysis steps



32

## Installing and loading R packages

**tidyverse**



`install.packages("tidyverse")`

Downloads files to computer

**1 x per computer**

`library("tidyverse")`

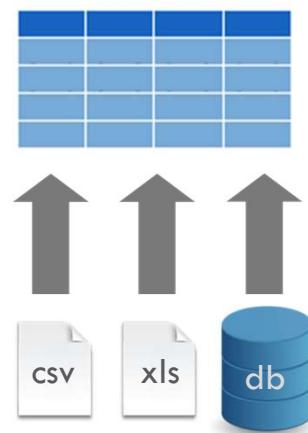
Loads package

**1 x per R Session**

33

## Dataframes: Beyond the Vector

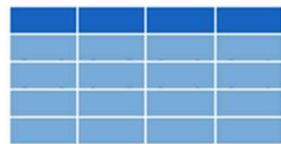
- Dataframe is the term for a table
- Dataframes are composed:  
Columns (Variables)  
Rows (Observations)
- Dataframes are objects and can be acted on like other objects



34

## read\_csv()

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

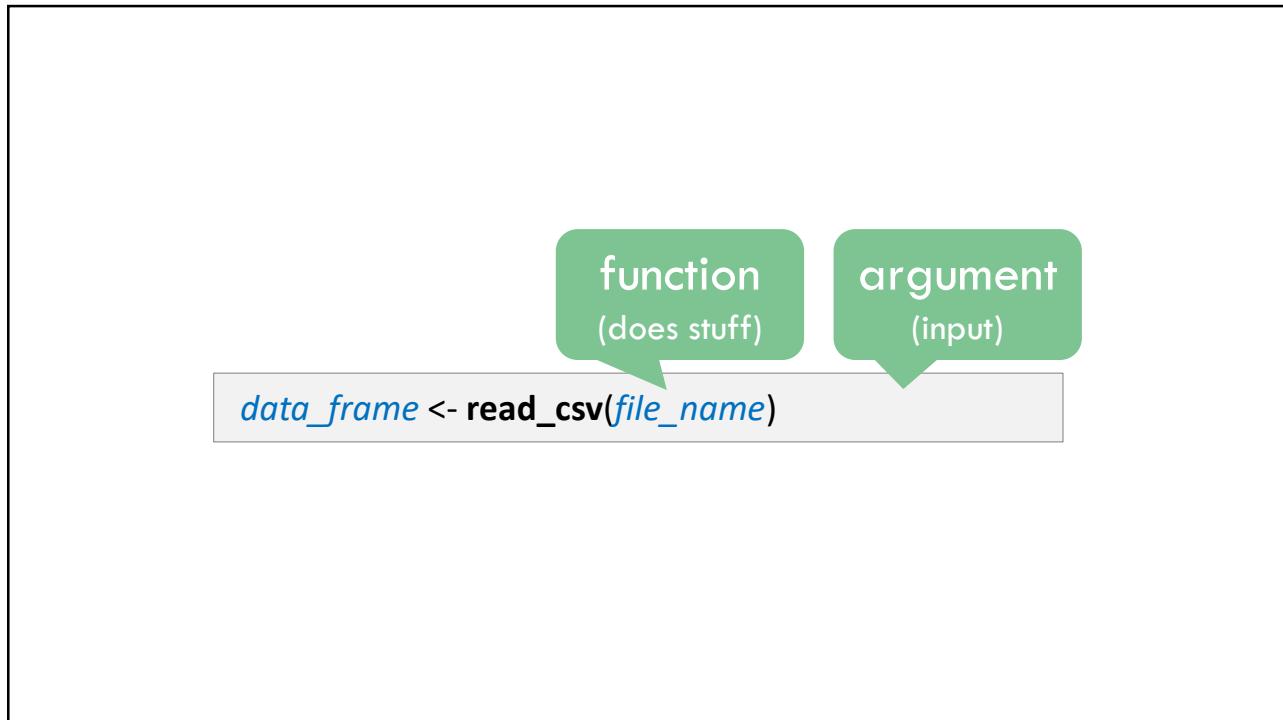


35

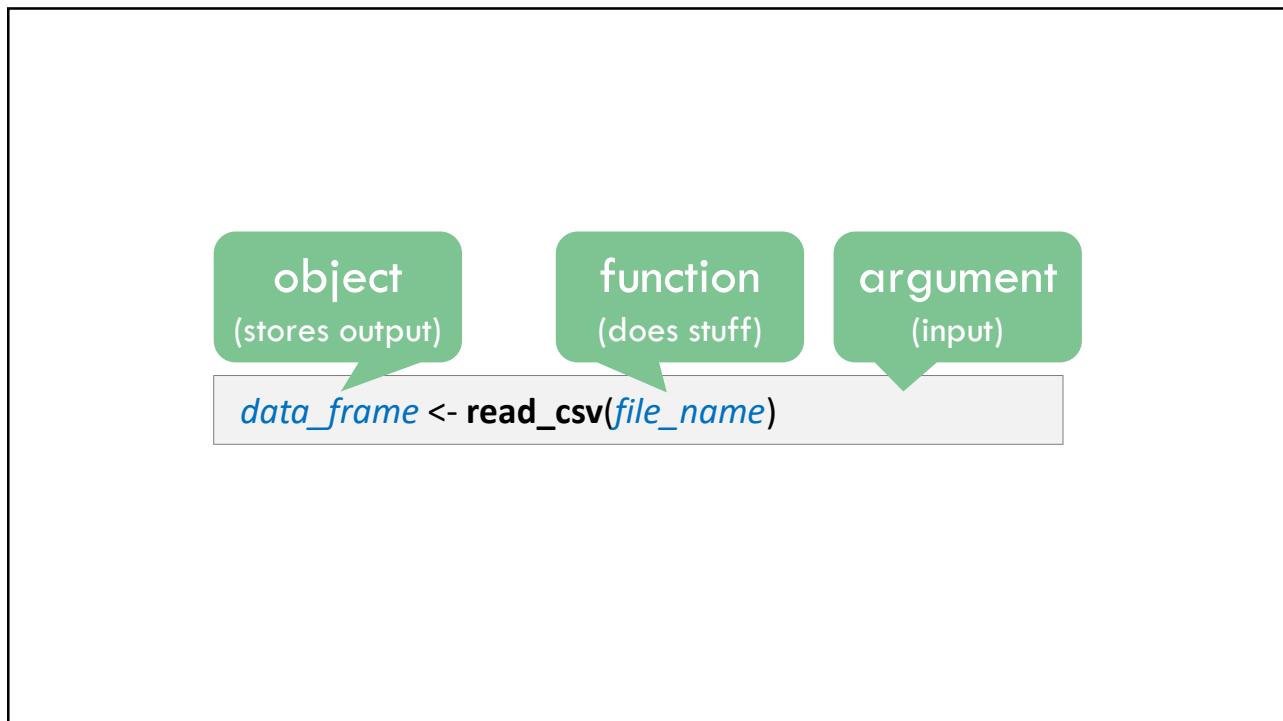
function  
(does stuff)

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

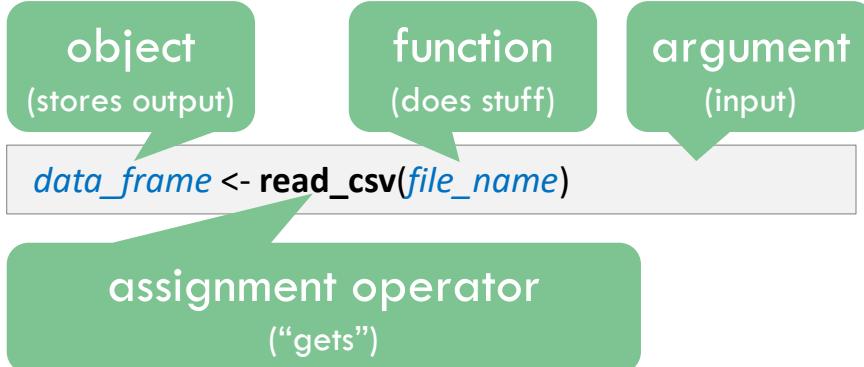
36



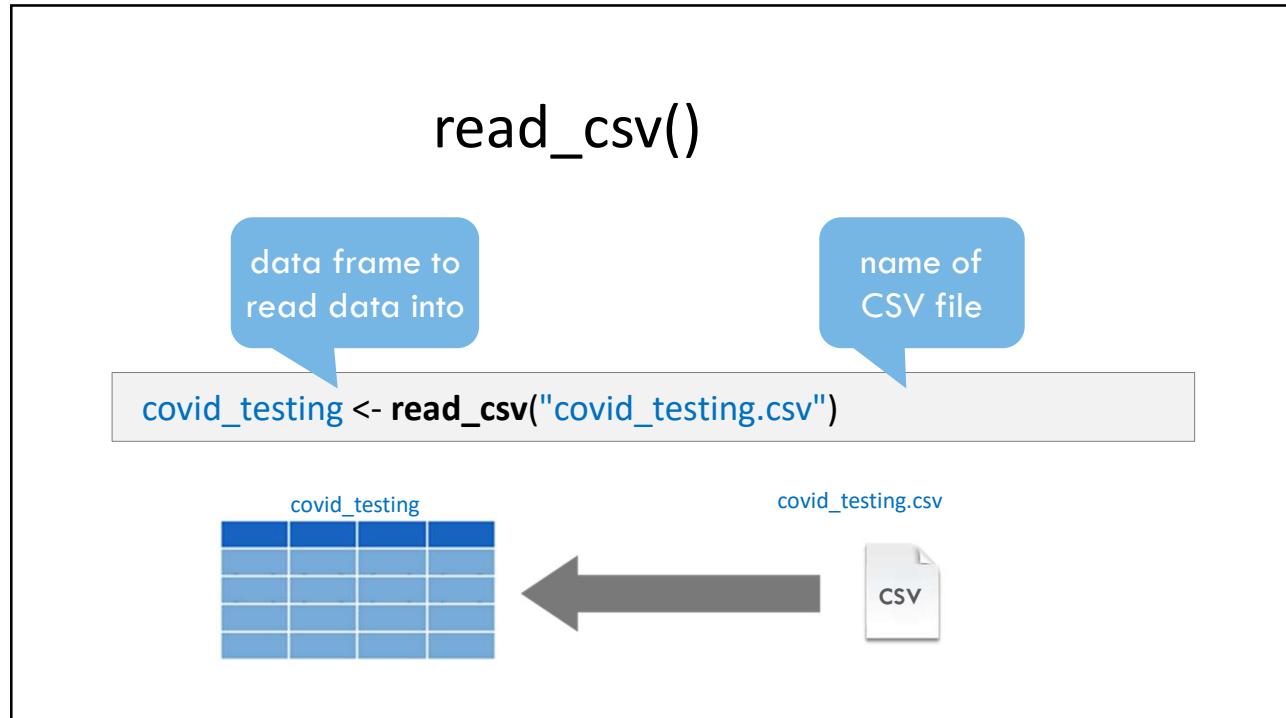
37



38



39



40

## Your Turn #4

In the MISC pane, select the folder:  
“exercises”

Select the R Markdown file:  
“01 - Importing and Exploring Data.Rmd”

In the Editor pane, follow the instructions to complete  
the exercise.

05 : 00

41

## Recap



Programming Language



IDE



Document Format

**Packages** extend the functionality of R. They need to be installed once per computer and loaded each session.

**Functions** do stuff. They accept **Arguments** to define parameters. We can store the output of functions in **Objects** using the assignment operator ( <- ).

**Importing Data** is the first step data analysis pipeline. `read_csv()` is a function from the tidyverse that we can use for importing data.

42

# What else?

43

## 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 **tidy** and to layout tidy data with **tidy**.

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)
- jsserialize** - json
- xml2** - XML
- httr** - Web APIs
- rvest** - HTML (Web Scraping)

### Save Data

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

**Common delimited file**  
`write_csv(x, path, na = "NA", append = FALSE, col_names = (append))`

**File with arbitrary delimiter**  
`write(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))`

### Save 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())`

**0.b.c** → **0, 1, 2, 3** Comma Delimited Files  
`read_csv("file.csv")`  
To mimic file.csv on  
`read_file(x = "file.csv", col1 = 1, col2 = 2, na = 5, NA, path = "file.csv")`

**0.b.c** → **0, 1, 2, 3** Semi-colon Delimited Files  
`read_csv("file.csv")`  
write\_file(x = "a;b;c\n1,2,3\n4,5,NA", path = "file.csv")

**a|b|c** → **0, 1, 2, 3** Fixed width  
`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")

**0.b.c** → **0, 1, 2, 3** 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")

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

**No header**  
`read_csv(file, col_names = FALSE)`

**Provide header**  
`read_csv(file, col_names = c("x", "y", "z"))`

**Read in a subset**  
`read_csv(file, n_max = 1)`

**Missing Values**  
`read_csv(file, na = c("1", ""))`

**Skip lines**  
`read_csv(file, skip = 1)`

**Read in a raw vector**  
`read_file(path)`

**Read each line into own vector**  
`read_file(path, skip = 0, max = 1L, na = character(), locale = default_locale(), progress = interactive())`

**Read lines into a file**  
`read_lines(file, skip = 0, n_max = 1L, progress = interactive())`

**Read Apache style log files**  
`read_logfile(file, col_names = FALSE, col_types = NULL, skip = 0, n_max = 1, progress = interactive())`

### Data types

**read** 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:

`#> #> col1 <- <class: integer>; #> age is an integer  
#> #> col2 <- <class: character>; #> sex is a character  
#> #> col3 <- <class: double>; #> bmi 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_ character()` the default
- `col_double()`, `col_double()`
- `col_datetime(format = "%Y-%m-%d")` Also `col_date(format = "%Y-%m-%d")`
- `col_date(format = "")`, `col_time(format = "")`
- `col_factor(levels, ordered = FALSE)`
- `col_logical()`
- `col_logical()`
- `col_numberic()`
- `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_double()` Also `parse_date()` and `parse_time()`
- `parse_double()`
- `parse_factor()`
- `parse_integer()`
- `parse_logical()`
- `parse_number()`

`x$A <- parse_number(x$A)`



44



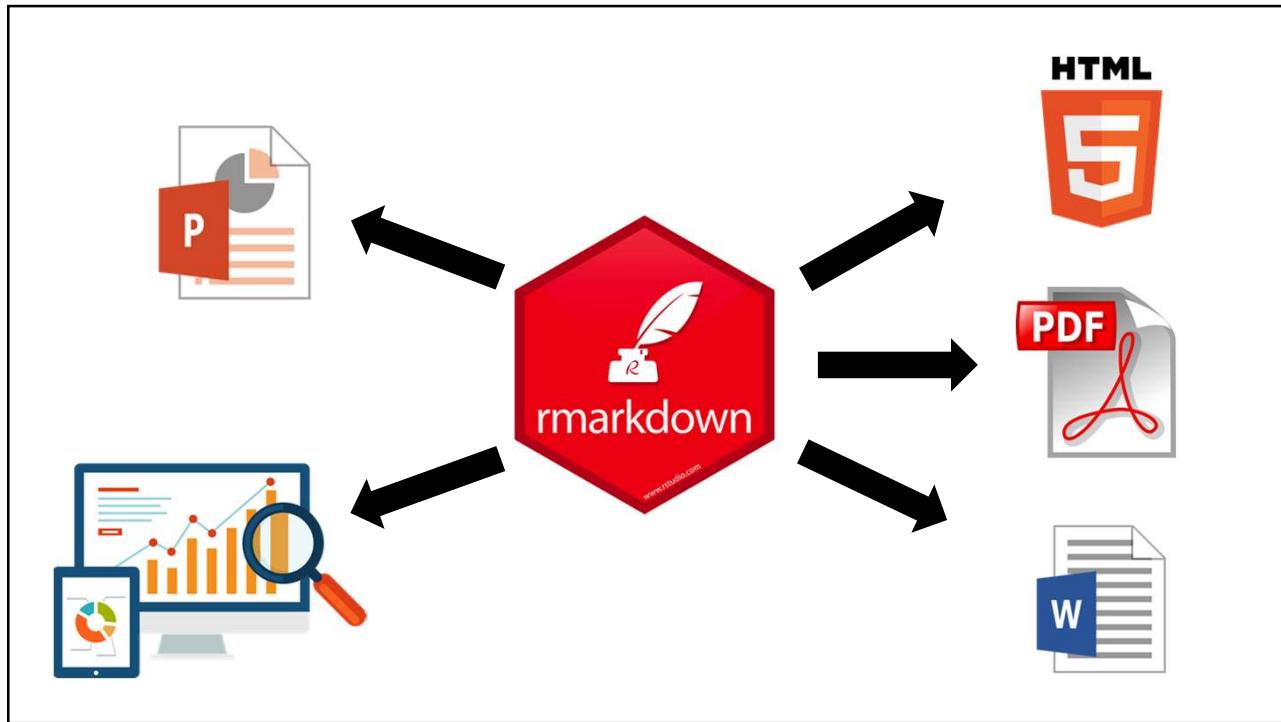
45

## Databases

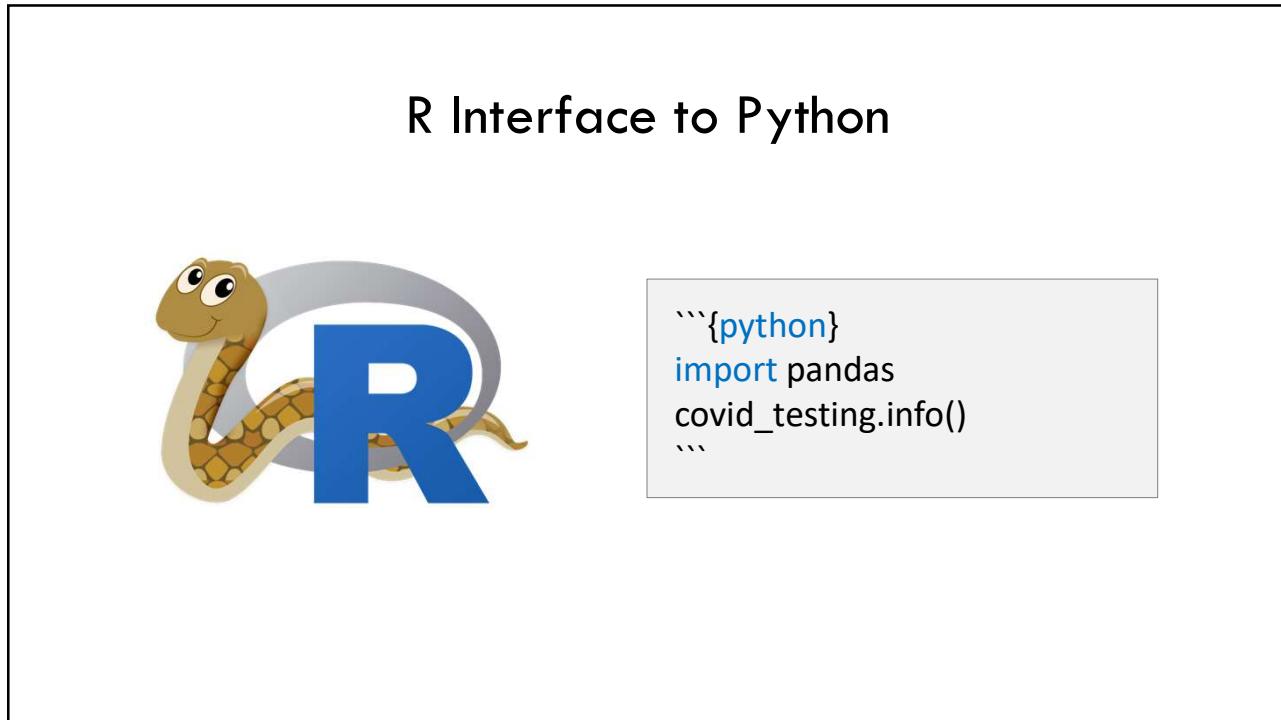
|   |                                      |   |                                 |   |                                 |
|---|--------------------------------------|---|---------------------------------|---|---------------------------------|
|  | <a href="#">Microsoft SQL Server</a> |  | <a href="#">Amazon Redshift</a> |  | <a href="#">Other Databases</a> |
|  | <a href="#">MonetDB</a>              |  | <a href="#">Apache Hive</a>     |  | <a href="#">PostgreSQL</a>      |
|  | <a href="#">MongoDB</a>              |  | <a href="#">Apache Impala</a>   |  | <a href="#">SQLite</a>          |
|  | <a href="#">MySQL</a>                |  | <a href="#">Athena</a>          |  | <a href="#">Salesforce</a>      |
|  | <a href="#">Netezza</a>              |  | <a href="#">Cassandra</a>       |  | <a href="#">Teradata</a>        |
|  | <a href="#">Oracle</a>               |  | <a href="#">Google BigQuery</a> |   |                                 |

<https://db.rstudio.com/databases/>

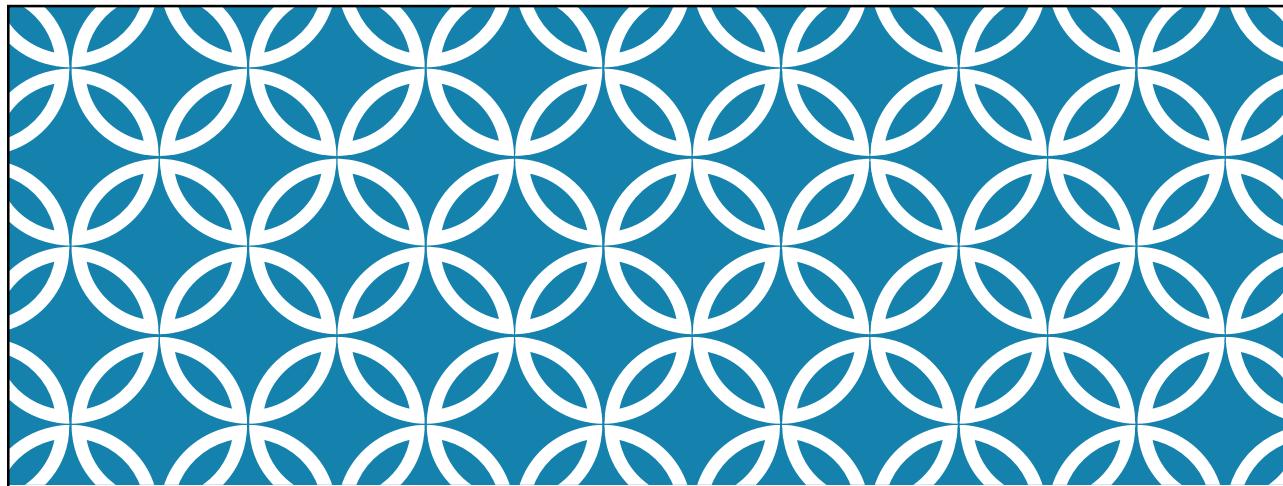
46



47



48



# Data Visualization

Session 3  
**Stephan Kadauke**

1

1

| Session  | Instructor      |
|--|-----------------|
| Instructor Introductions, Introduction to technology | Amrom Obstfeld  |
| Introduction to R and RStudio                        | Joe Rudolf      |
| Reproducible Reporting                               | Joe Rudolf      |
| Data Visualization                                   | Stephan Kadauke |
| Data Transformation                                  | Amrom Obstfeld  |
| Statistical Analysis                                 | Dan Herman      |
| Advanced Reporting                                   | Patrick Mathias |

2

2

## Goals

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

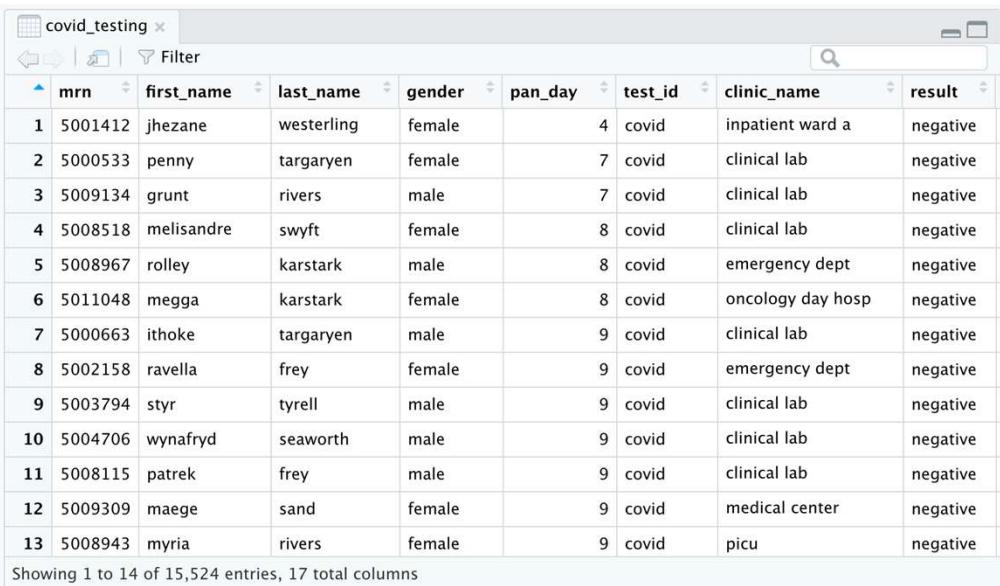
## Objectives

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

3

3

## covid\_testing



| mrn | first_name | last_name  | gender    | pan_day | test_id | clinic_name      | result            |          |
|-----|------------|------------|-----------|---------|---------|------------------|-------------------|----------|
| 1   | jhezane    | westerling | female    | 4       | covid   | inpatient ward a | negative          |          |
| 2   | penny      | targaryen  | female    | 7       | covid   | clinical lab     | negative          |          |
| 3   | grunt      | rivers     | male      | 7       | covid   | clinical lab     | negative          |          |
| 4   | 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

4

4

## Your Turn 1

Consider the `covid_testing` data frame.

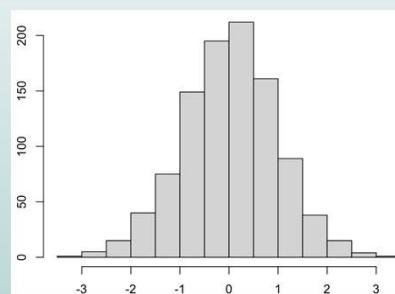
What do you think plot would look like in which:

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

5

5

## Your Turn 2



What is the name of this kind of plot?

Type the answer into the chat!

6

6

## 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))
```

7

7

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

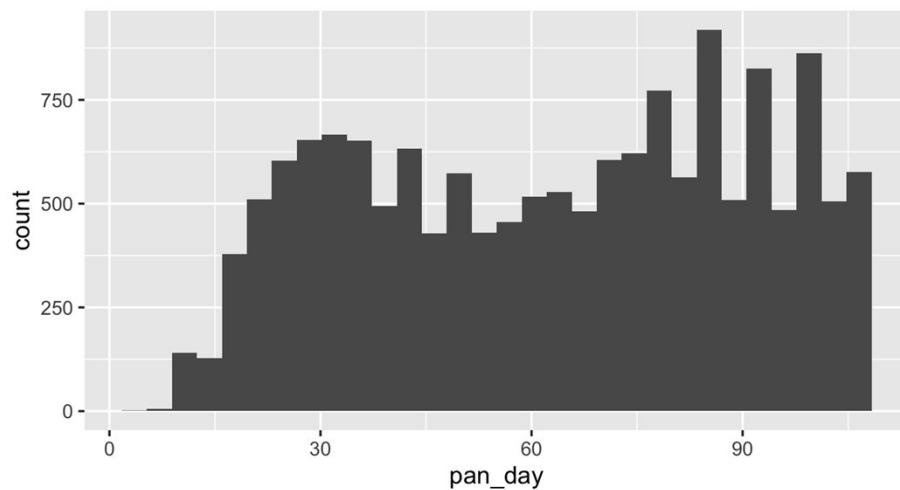
Often times, things that look like an error in R are actually just a message.

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.

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

8

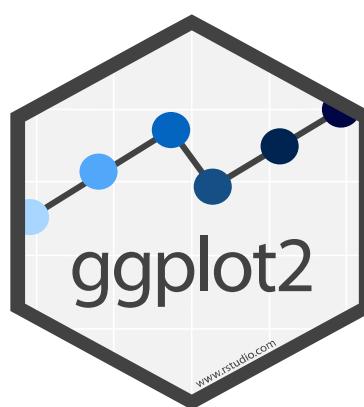
8



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

9

9



10

10

# 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

11

11

## To make **any** kind of graph:

1. Pick a “tidy”  
data frame

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

2. Pick a “geom”  
function

3. Write aesthetic  
mappings

12

1. Pick a “tidy”  
data frame

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

13

13

## 1. Pick a “Tidy” Data Frame

A data set is **tidy** if:

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

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

Wickham H. **Tidy Data**. J Stat Soft 2014.

14

14

1. Pick a “tidy”  
**data frame**

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

2. Pick a “geom”  
function

15

## 2. Pick a “Geom” Function



`geom_histogram()`



`geom_dotplot()`



`geom_bar()`



`geom_boxplot()`



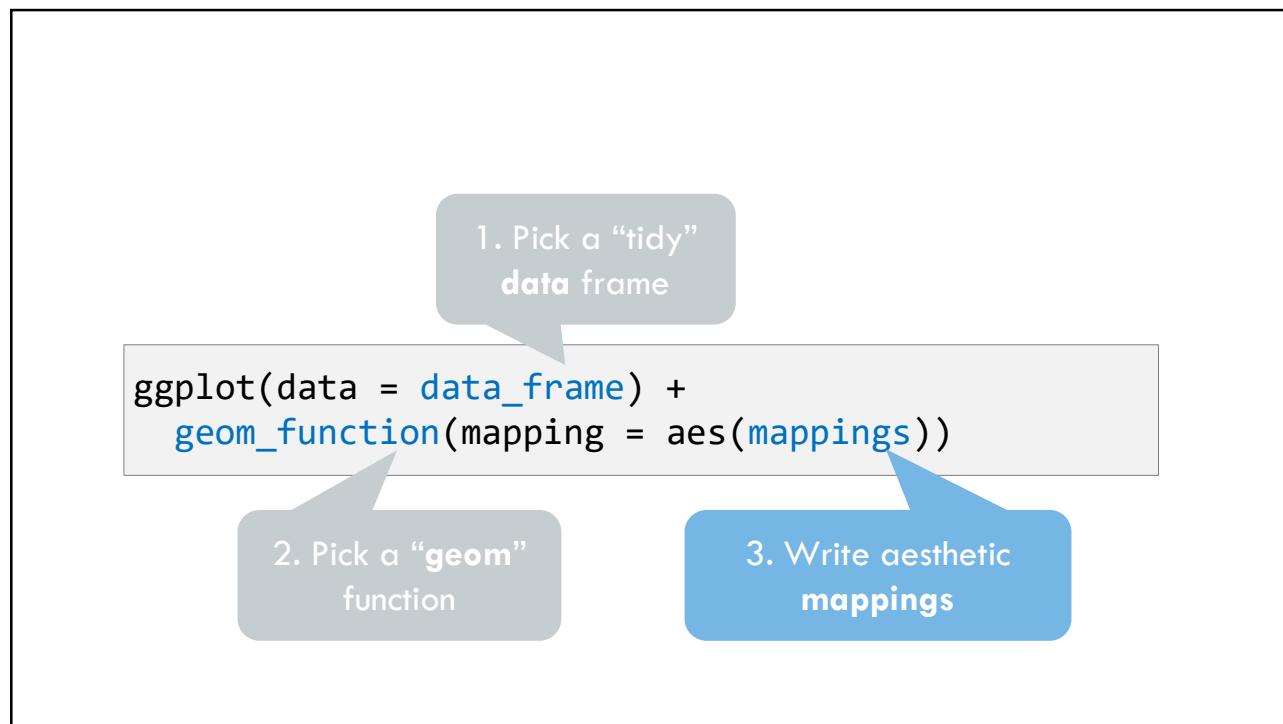
`geom_point()`



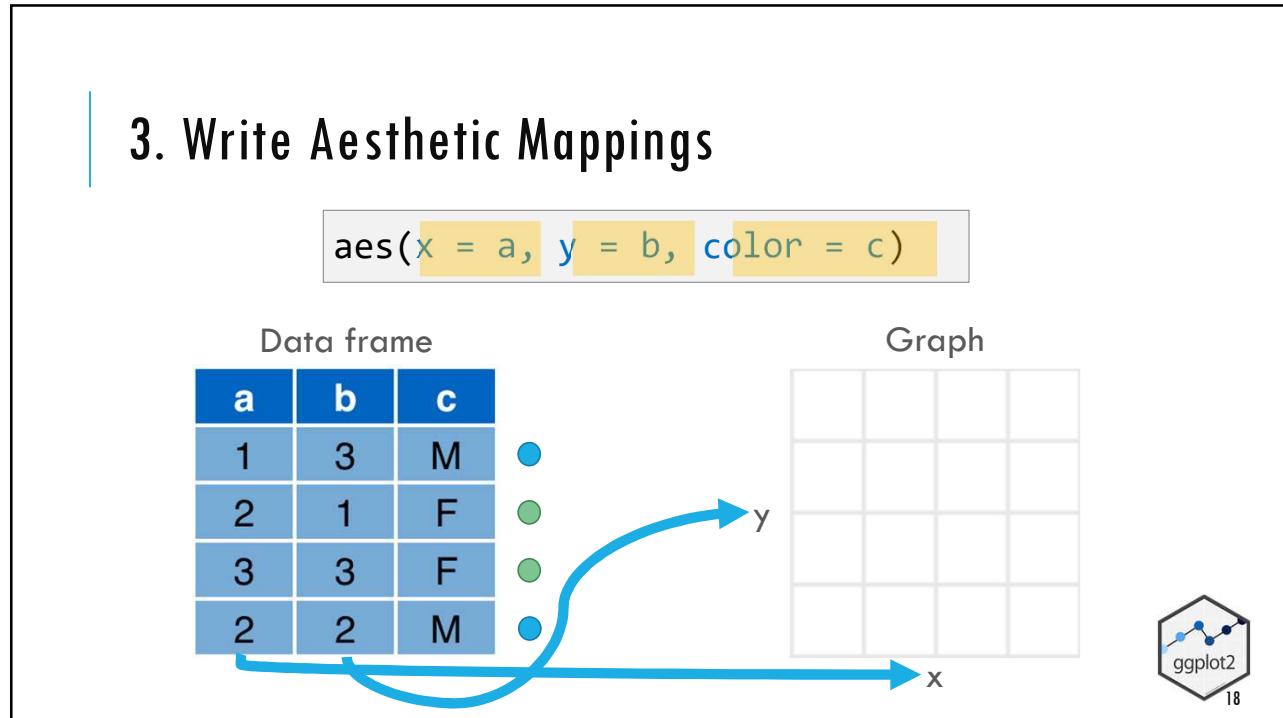
`geom_line()`

16

16

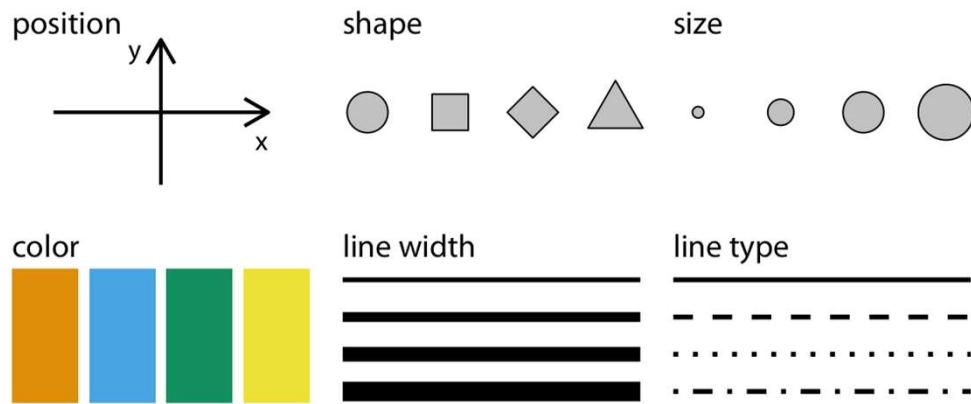


17



18

## Aesthetics



From [Fundamentals of Data Visualization](#), by Claus Wilke, licensed under CC-BY-NC-ND

19

19

## 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



20

10

## Your Turn 4

Open 03 - Visualize.Rmd. Work through the exercises of the section titled “Your Turn 4”.

Stop when it says “Stop Here”.

Click “yes” when you’re done!

21

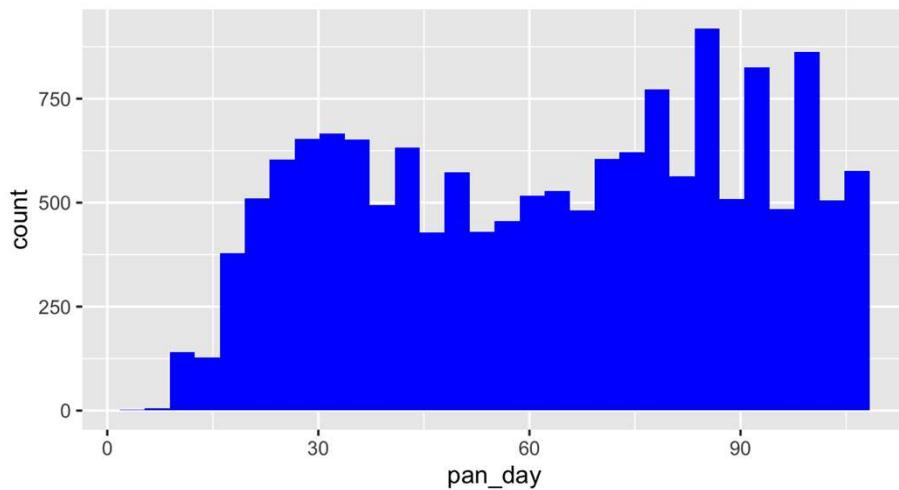
21

## Setting vs Mapping Aesthetics

22

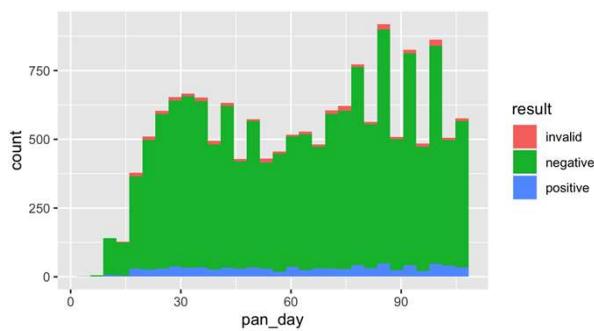
11

How would you make this plot?



23

23

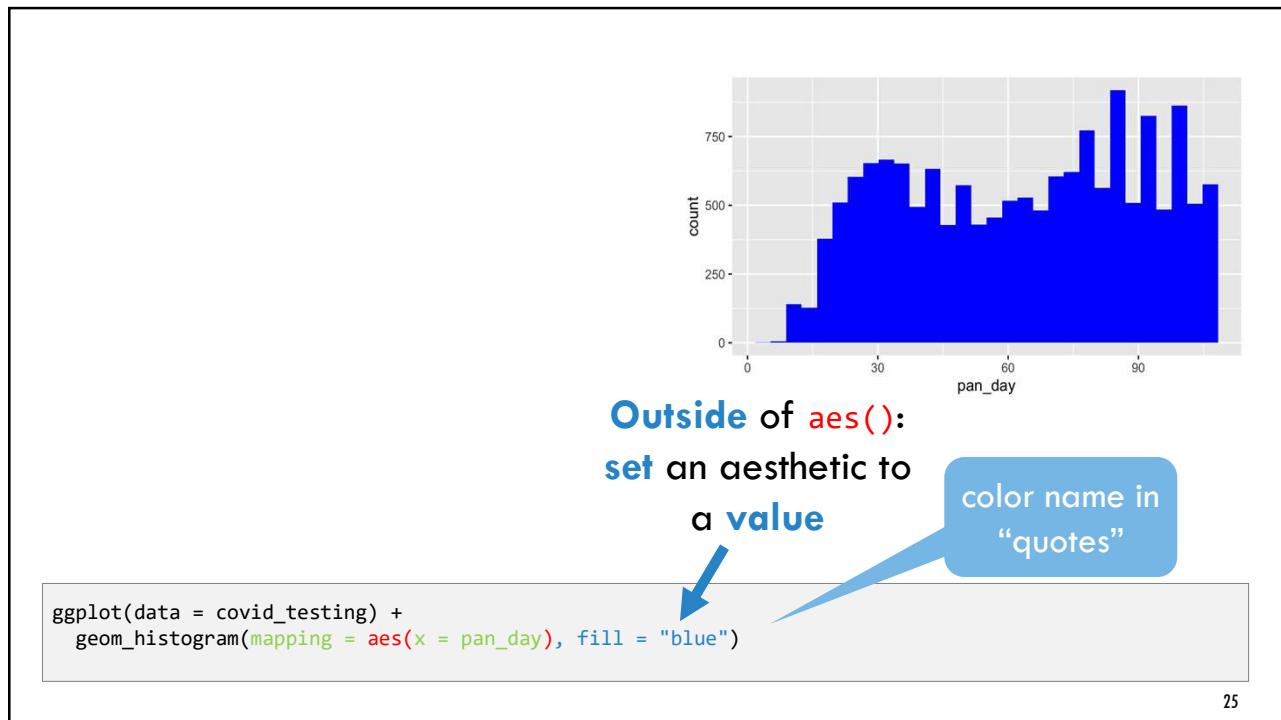


Inside of `aes()`:  
map an aesthetic  
to a variable

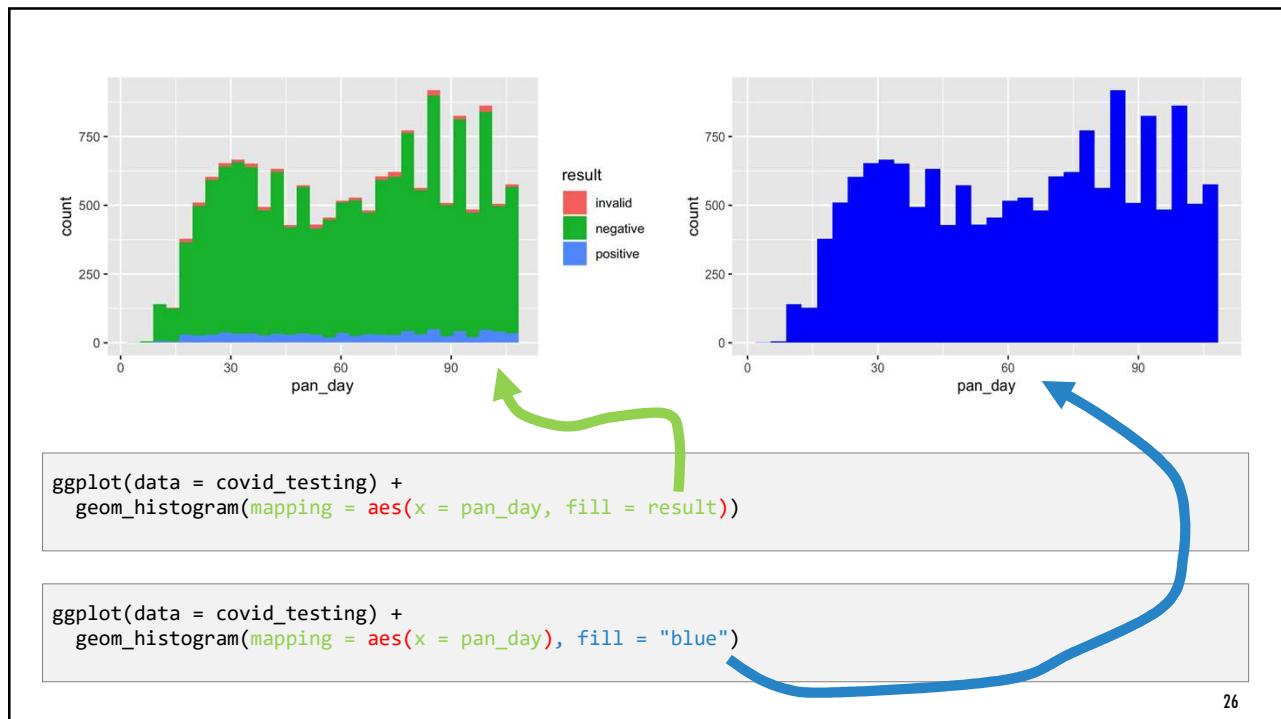
```
ggplot(data = covid_testing) +  
  geom_histogram(mapping = aes(x = pan_day, fill = result))
```

24

24



25

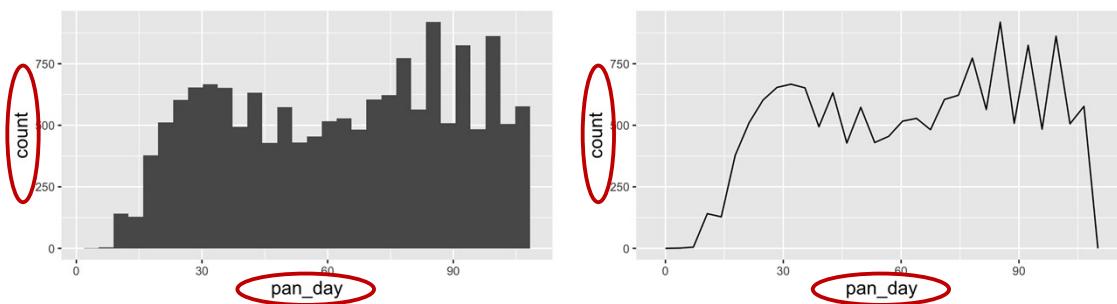


26

# Geom Functions

27

How are these plots similar?

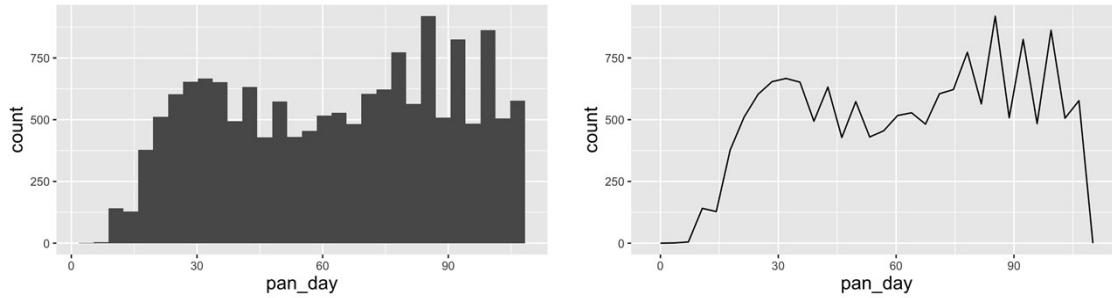


Same: x axis, y axis, data

28

28

How are these plots different?



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

29

29

## Your Turn 5

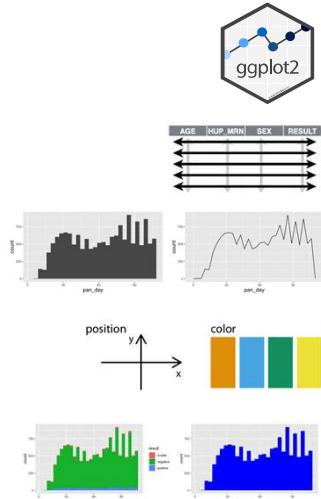
Return to [03-Visualize.Rmd](#). Work through the exercises of the section titled “Your Turn 5.”

Click “yes” when you’re done!

30

30

## Recap



**ggplot2** is a package that provides a **grammar of graphics**. You can create **any type of plot** using a simple template to which you provide:

1. A **tidy data frame**, in which each **variable** is in its own **column**, each **observation** is in its own **row**, each **value** is in its own **cell**;
2. A **geom function**, which tells R what kind of plot to make; and
3. **Aesthetic mappings**, which tell R how to represent data as graphical markings on the plot.

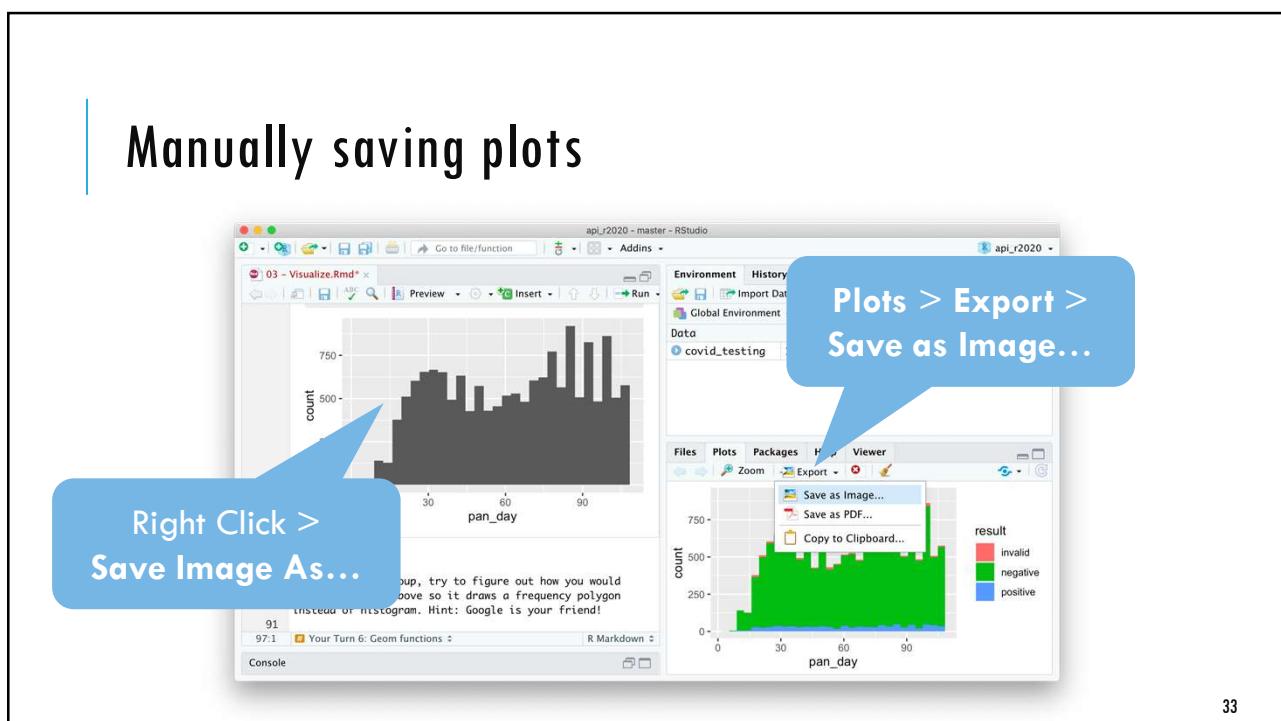
Aesthetics can be **mapped** to a variable or **set** to a constant value.

31

What Else?

32

## Manually saving plots

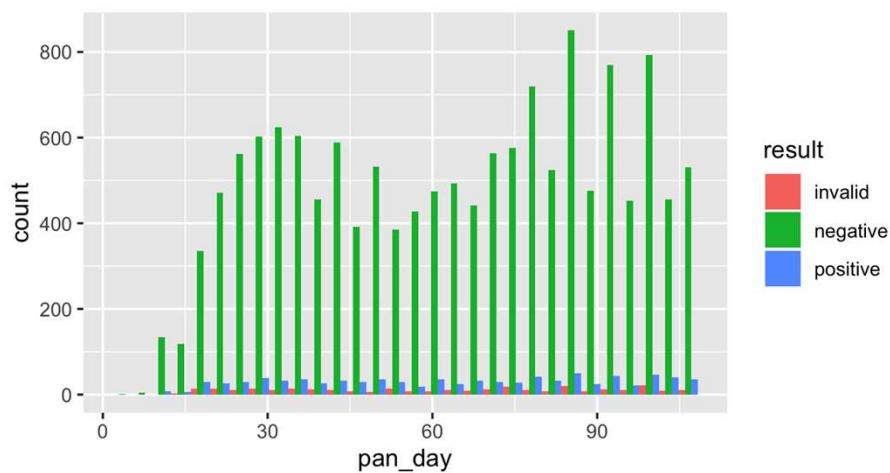


33

33

## Position adjustments

How overlapping objects are arranged

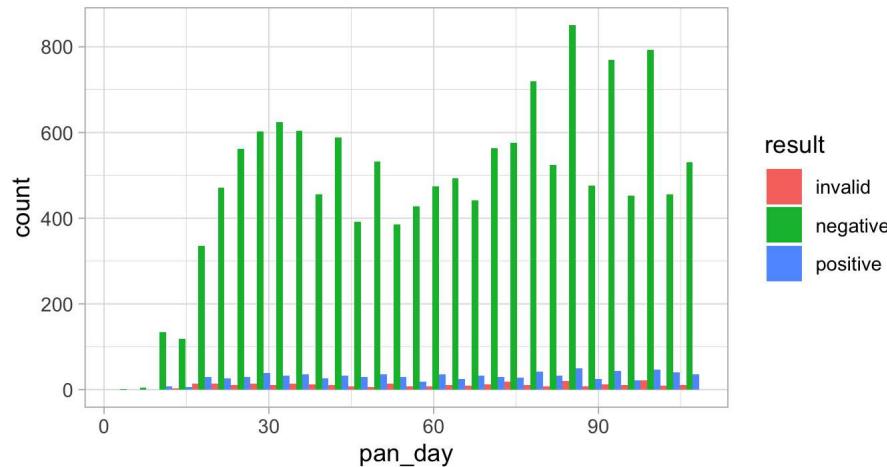


34

34

# Themes

Visual appearance of non-data elements

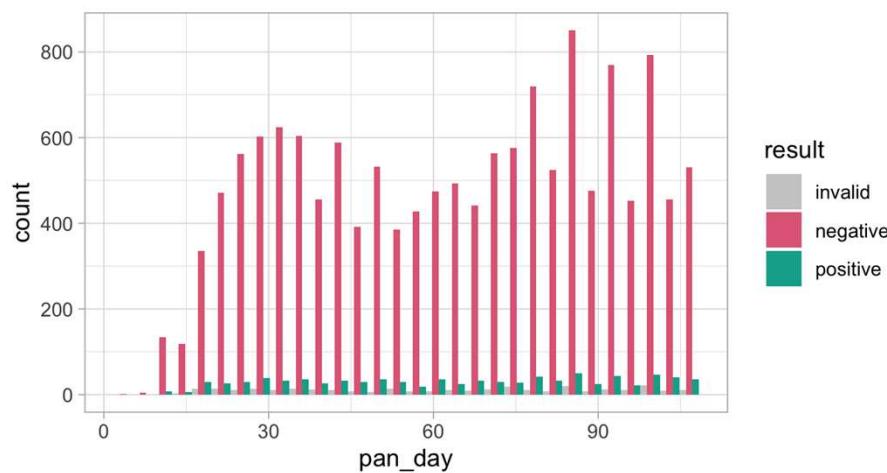


35

35

# Scales

Customize color scales and other mappings

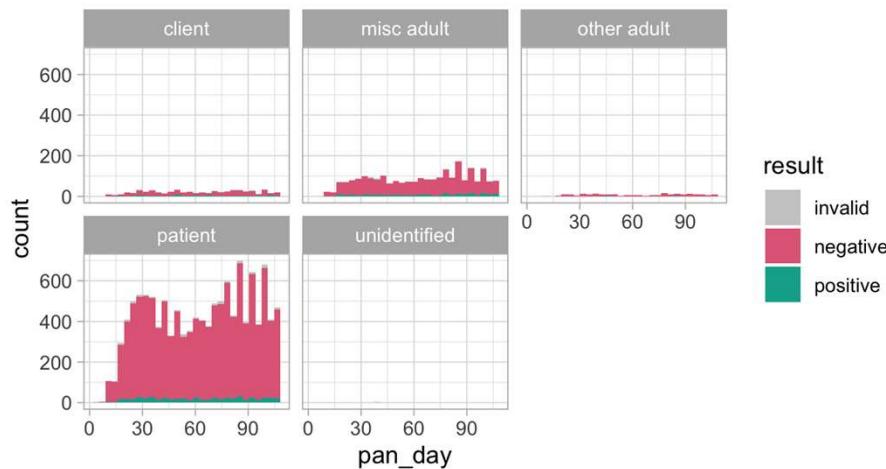


36

36

# Facets

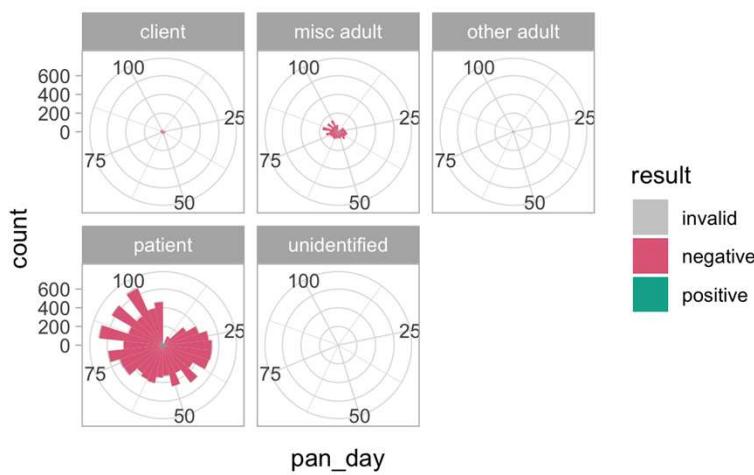
Subplots that display subsets of the data



37

37

# Coordinate systems



38

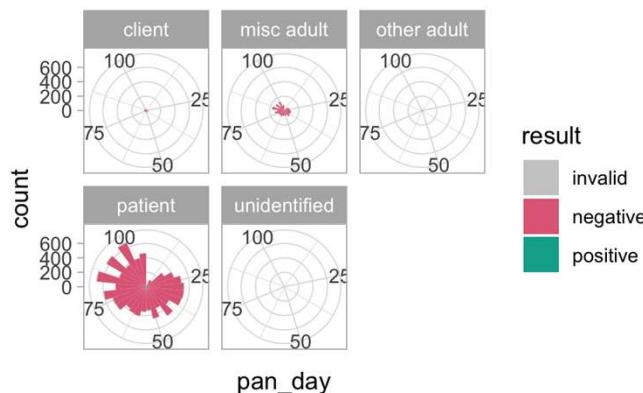
38

19

# Titles and captions

COVID19 test volume

Displayed in polar coordinates, mostly to show off ggplot2



39

39

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

Required

Optional

40

40

<https://r4ds.had.co.nz/>

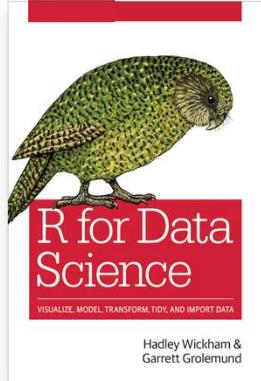
R for Data Science

Table of contents

- Welcome**
- 1 Introduction
- Explore
- 2 Introduction
- 3 Data visualisation
- 4 Workflow: basics
- 5 Data transformation
- 6 Workflow: scripts
- 7 Exploratory Data Analysis
- 8 Workflow: projects
- Wrangle
- 9 Introduction
- 10 Tibbles
- 11 Data import

## Welcome

This is the website for “R for Data Science”. This book will teach you how to do data science with R: You’ll learn how to get your data into R, get it into the most useful structure, transform it, visualise it and model it. In this book, you will find a practicum of skills for data science. Just as a chemist learns how to clean test tubes and stock a lab, you’ll learn how to clean data and draw plots—and many other things besides. These are the skills that allow data science to happen, and here you will find the best practices for doing each of these things with R. You’ll learn how to use the grammar of graphics, literate programming, and reproducible research to save time. You’ll also learn how to manage cognitive resources to facilitate discoveries when wrangling, visualising, and exploring data.



41

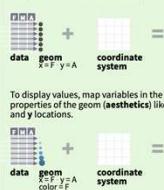
41

<https://www.rstudio.com/resources/cheatsheets/>

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.



Complete the template below to build a graph.

```
ggplot(data = <DATA>) +
  <GEO FUNCTION>(mapping = aes(<MAPPINGS>,
  stat = <STAT>, position = <POSITION>)) +
  <COORDINATE FUNCTION> +
  <FACTORY FUNCTION> +
  <SCALE FUNCTION> +
  <THEME FUNCTION>
```

Required

Not required, requires defaults supplied

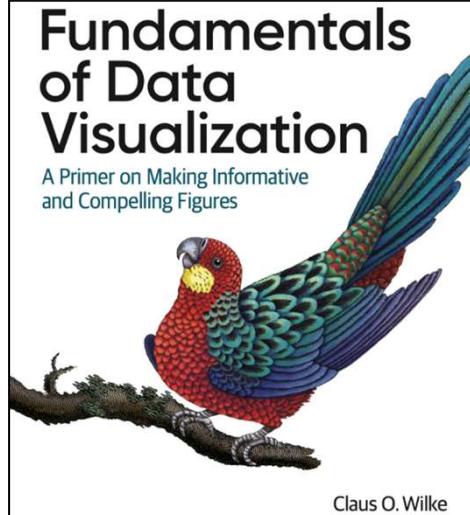
ggplot(data = mpg, aes(x = cyl, y = hwy)) Begins a plot

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

|  |   |
|--|---|
| <b>GRAPHICAL PRIMITIVES</b>  | <b>TWO VARIABLES</b>  |
| a <- ggplot(teens, aes(x=age, y=unemploy))<br>b <- ggplot(teens, aes(x=long, y=lat))<br><br>a + geom_point()<br>(Useful for expanding limits)<br>b + geom_rect(yend=lat+1, curvature=1) -> x, xend, y, yend,<br>alpha, angle, curve, linetype, size<br>a + geom_path(needle="butt", linejoin="round",<br>linecap="square")<br>x, y, alpha, color, group, linetype, size<br>a + geom_polygon(segments=group = group)<br>x, y, alpha, color, fill, group, linetype, size<br>b + geom_rect(rectxmin = long, ymin=lat, xmax=<br>long + 1, ymax = lat + 1) -> x, xmax, xmin, ymax,<br>ymin, alpha, color, fill, linetype, size<br>a + geom_ribbon(ymin=unemploy - 900,<br>ymax=unemploy + 900) -> x, ymax, ymin,<br>alpha, color, fill, group, linetype, size | continuous x, continuous y<br>e <- ggplot(economics, aes(date, unemploy))<br><br>e + geom_label(label = "cty", nudge_x = 1,<br>nudge_y = 1, check_overlap = TRUE) x, y, label,<br>alpha, angle, color, family, fontface, hjust,<br>lineheight, size<br>e + geom_interpolate(x, y, alpha, color, fill, shape,<br>size, stroke, strokewidth, width)<br>e + geom_point(x, y, alpha, color, fill, shape,<br>size, stroke, strokewidth, width)<br>e + geom_quantile(x, y, alpha, color, group,<br>linetype, size, weight)<br>e + geom_rug(sides = "bl") x, y, alpha, color,<br>linetype, size<br>e + geom_smooth(method = lm) x, y, alpha,<br>color, fill, group, linetype, size, weight<br>e + geom_text(label = "cty", nudge_x = 1,<br>nudge_y = 1, check_overlap = TRUE) x, y, label,<br>alpha, angle, color, family, fontface, hjust,<br>lineheight, size, vjust<br><br>discrete x, continuous y<br>f <- ggplot(mpg, aes(class, hwy))<br><br>f + geom_bar(stat = "bin")<br>x, y, alpha, color, fill, linetype, size<br>c <- ggplot(mpg, aes(hwy))<br>c + geom_area(stat = "bin")<br>x, y, alpha, color, fill, linetype, size |
| LINE SEGMENTS  | continuous bivariate distribution   |
| common aesthetics: x, y, alpha, color, linetype, size  | h <- ggplot(diamonds, aes(carat, price))<br><br>h + geom_bin2d(binwidth = c(0.25, 500))<br>x, y, alpha, color, fill, linetype, size<br>h + geom_density2d()<br>x, y, alpha, color, fill, linetype, size<br>h + geom_hex()<br>x, y, alpha, color, fill, linetype, size   |
| b + geom_abline(intercept=0, slope=1)<br>b + geom_hline(intercept=lat)<br>b + geom_vline(intercept=long)   | continuous function   |
| b + geom_segment(aes(endx=lat+1, xend=long+1))<br>b + geom_spline()  | i <- ggplot(economics, aes(date, unemploy))<br><br>i + geom_area()<br>x, y, alpha, color, fill, linetype, size<br>i + geom_line()<br>x, y, alpha, color, group, linetype, size<br>i + geom_step(direction = "hv")<br>x, y, alpha, color, group, linetype, size  |
| ONE VARIABLE continuous  | visualizing error   |
| c <- ggplot(mpg, aes(hwy))<br>c2 <- ggplot(mpg)  | j <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)<br>j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))<br><br>j + geom_crossbar(fatten = 2)<br>x, y, alpha, color, fill, group, linetype, size<br>j + geom_errorbar()<br>x, y, alpha, color, fill, group, linetype, size<br>j + geom_linerange()<br>x, y, alpha, color, fill, group, linetype, size<br>j + geom_pointrange()<br>x, y, alpha, color, fill, group, linetype, size  |

42

42



<https://clauswilke.com/dataviz/>

43

43

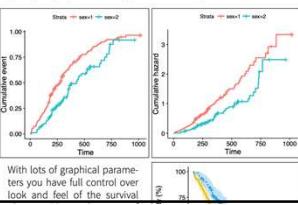
## Creating Survival Plots Informative and Elegant with survminer

### Survival Curves

```
library("survival")
fit <- survfit(Surv(time,status) ~ sex, data = lung)
class(fit)
## [1] "survfit"
library("survminer")
ggsurvplot(fit, data = lung)
```

Use the `fun` argument to set the transformation of the survival curve. E.g. `"event"` for cumulative events, `"cumhaz"` for the cumulative hazard function or `"pct"` for survival probability in percentage.

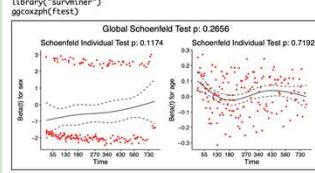
```
ggsurvplot(fit, data = lung, fun = "event")
ggsurvplot(fit, data = lung, fun = "cumhaz")
```



### Diagnostics of Cox Model

The function `cox.zph()` from `survival` package may be used to test the proportional hazards assumption for a Cox regression model fit. The graphical verification of this assumption may be performed with the function `ggoxph()` from the `survminer` package. For each covariate it produces plots with scaled Schoenfeld residuals against the time.

```
library("survival")
fit <- coxph(Surv(time, status) ~ sex+age, data = lung)
ftest
library("survminer")
ggoxph(fit)
```



The function `ggcoxdiagnostics()` plots different types of residuals as a function of time, linear predictor or observation id. The type of residual is selected with `type` argument. Possible values are `"martingale"`, `"deviance"`, `"score"`, `"schoenfeld"`, `"dbeta"`, `"dfbetas"`, and `"scaledsch"`. The `ox.scale` argument defines what shall be plotted on the X axis. Possible values are `"linear.predictions"`, `"observation.id"`, `"time"`.

Logical arguments `hline` and `sline` may be used to add horizontal line or smooth line to the plot.

```
library("survival")
library("survminer")
fit <- coxph(Surv(time, status) ~ sex + age, data = lung)

ggcoxdiagnostics(fit,
  type = "deviance",
  ox.scale = "linear.predictions")
ggcoxdiagnostics(fit,
```

### Summary of Cox Model

The function `ggforest()` from the `survminer` package creates a forest plot for a Cox regression model fit. Hazard ratio estimates along with confidence intervals and p-values are plotted for each variable.

```
library("survival")
library("survminer")
lungstage <- ifelse(lung$stage > 70, ">70", "<=70")
lung <- coxph(Surv(time, status) ~ sex + ph.ecog + age, data = lung)
fit
```

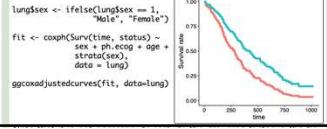
```
# Call:
#> coxph(formula = Surv(time, status) ~ sex+ph.ecog+age, data=lung)
#>   coef exp(coef) se(coef)    z      p
#> sex     -0.567    0.567   0.165 -3.37 0.00075
#> ph.ecog  0.307    1.399   0.187  1.64 0.10375
#>
#> Likelihood ratio test=31.6 on
#>   227 number of events= 164
```



The function `ggadjustedcurves()` from the `survminer` package plots Adjusted Survival Curves for Cox Proportional Hazards Model. Adjusted Survival Curves show how a selected factor influences survival estimated from the Cox model.

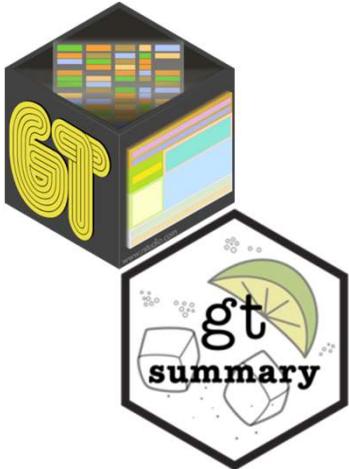
```
library("survival")
library("survminer")
```

```
lungsex <- ifelse(lung$sex == 1,
  "Male", "Female")
fit <- coxph(Surv(time, status) ~ sex + ph.ecog + age +
  (1|lung$id), data = lung)
ggadjustedcurves(fit, data=lung)
```



44

## A grammar for tables



| Variable              | N   | Drug A, N = 98 <sup>1</sup> | Drug B, N = 102 <sup>1</sup> | p-value <sup>2</sup> |
|-----------------------|-----|-----------------------------|------------------------------|----------------------|
| <b>Age</b>            | 189 | 46 (37, 59)                 | 48 (39, 56)                  | 0.7                  |
| <b>Grade</b>          | 200 |                             |                              | 0.9                  |
| I                     |     | 35 (36%)                    | 33 (32%)                     |                      |
| II                    |     | 32 (33%)                    | 36 (35%)                     |                      |
| III                   |     | 31 (32%)                    | 33 (32%)                     |                      |
| <b>Tumor Response</b> | 193 | 28 (29%)                    | 33 (34%)                     | 0.6                  |

<sup>1</sup> Statistics presented: median (IQR); n (%)

<sup>2</sup> Statistical tests performed: Wilcoxon rank-sum test; chi-square test of independence

45

## Goals

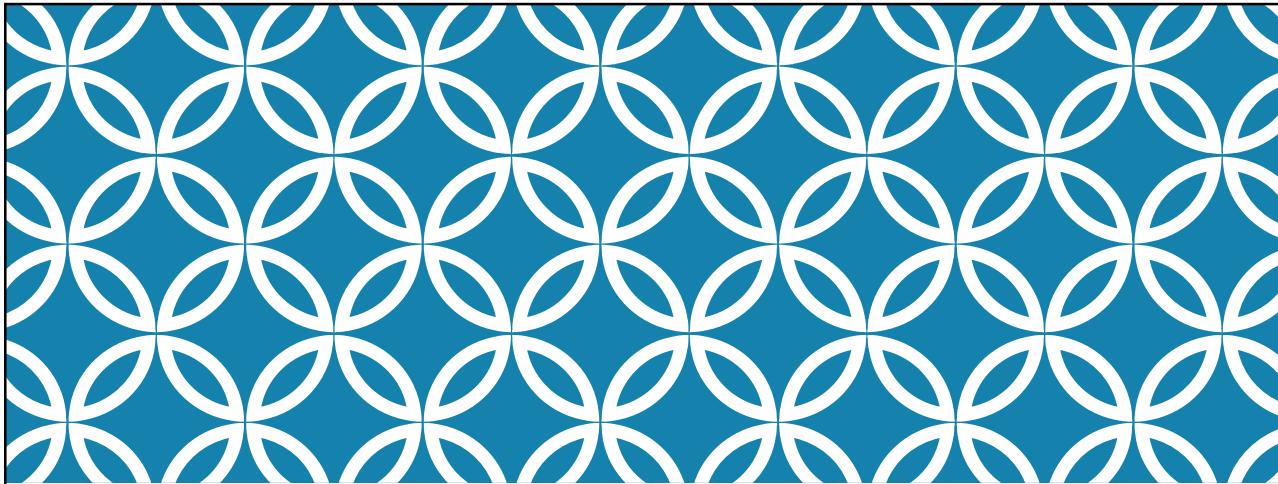
1. Appreciate the importance of visualization for understanding data
2. Learn how to use ggplot2 to visualize data

## 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

46

46



# Data Transformation

Session 4  
Amrom Obstfeld

1

## Goal

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

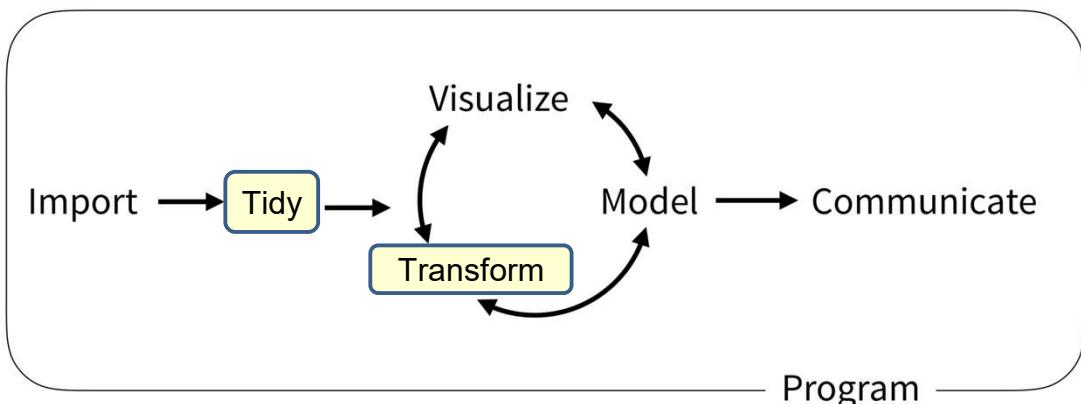
## Objectives

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

2

2

# Typical Data Science Pipeline



3

## What is a “Tidy” Data Frame

A data set is **tidy** if:

| AGE | MRN | SEX | RESULT |
|-----|-----|-----|--------|
|     |     |     |        |
|     |     |     |        |
|     |     |     |        |
|     |     |     |        |

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

4

4

## 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("covid_testing.csv")
```

5

## Pop Quiz

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

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

6

6

# Transform Data with



7

7

# dplyr



dplyr implements a *grammar* for transforming tabular data.



8

8

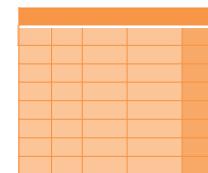
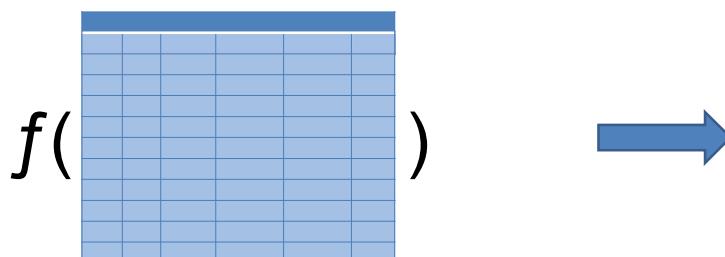
## dplyr: a grammar for transforming data

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



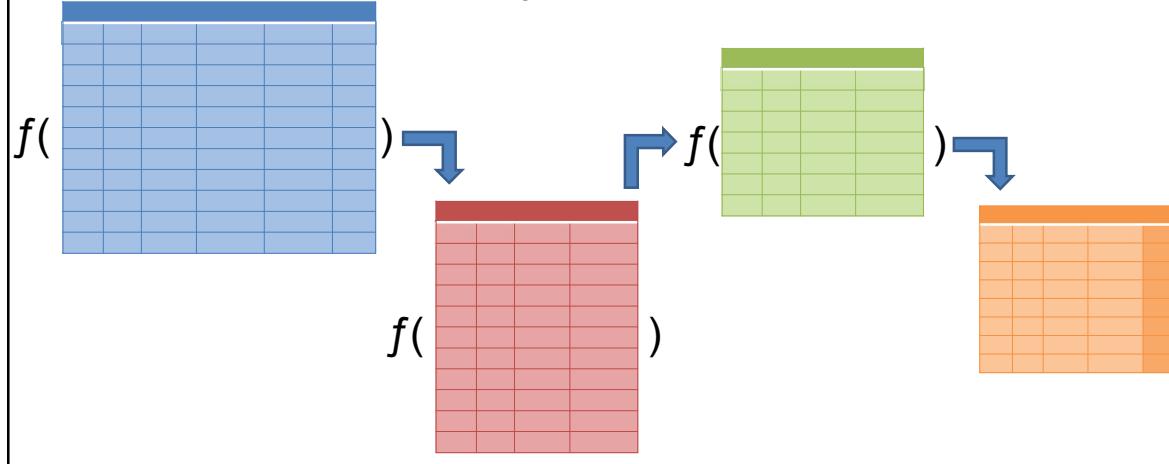
9

## Dplyr Approach



10

## Dplyr Approach



11

11

## 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



12

12

## Isolating data

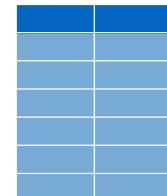
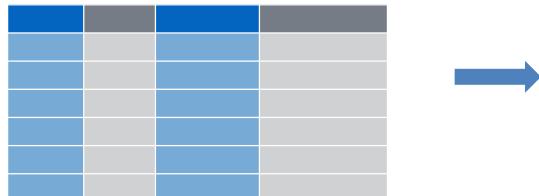
13

select()

14

# select()

Extract columns from a data frame



= Number of rows  
↓ Number of Columns



15

15

# 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)



16

16

## 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  |



17

17

## 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 |

...

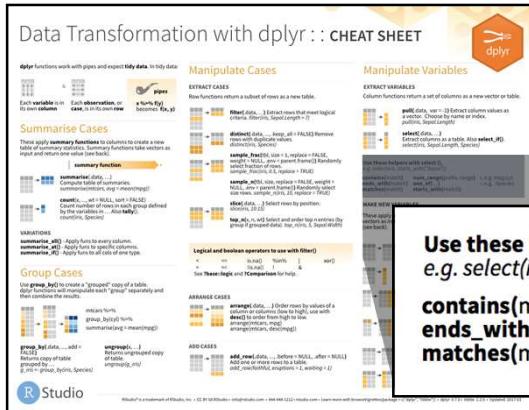
| first_name | gender |
|------------|--------|
| sarella    | female |
| alester    | male   |
| jhezane    | female |
| penny      | female |



18

18

# select() helpers



The image shows a 'Data Transformation with dplyr :: CHEAT SHEET'. It is a reference card for dplyr functions, organized into sections: Summarise Cases, Group Cases, Manipulate Cases, and Add Cases. Each section contains icons and brief descriptions of the functions. A large callout box highlights the 'select()' helper functions:

**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)**

19



## Your Turn 2

- Alter the code to select just the `first_name` column from `covid_testing`
- Use the second code chunk to see if you can remove the `first_name` column

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

20

20

# filter()

21

## filter()

Extract rows that meet logical criteria

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



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

↓ Number of rows  
= Number of Columns



22

22

# 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



23

23

## 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 |



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



24

24

# 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 |



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



25

25

# 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     | cleagane   |



|  | mrn     | first_name | last_name |
|--|---------|------------|-----------|
|  | 5000083 | lollys     | cleagane  |



26

26

# 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)



27

27

# 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     |



28

28

# 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)



29

29

## Logical tests

|                        |                          |
|------------------------|--------------------------|
| <code>x &lt; y</code>  | Less than                |
| <code>x &gt; y</code>  | Greater than             |
| <code>x == y</code>    | Equal to                 |
| <code>x &lt;= y</code> | Less than or equal to    |
| <code>x &gt;= 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                |



30

30

# Pop Quiz

What is the result?

`1 == 1`

31

31

# Pop Quiz

What is the result?

`3 != 1`

32

32

# filter() variants

**Data Transformation with dplyr :: CHEAT SHEET**

**Manipulate Cases**

- Summarise Cases**: Tidy syntax summary functions to columns to create a new table. Use `summarise(x, ...)` or `summarise_(x, ..., .by_group = TRUE)`.
- VARIATIONS**: `summarise_all()`, `summarise_if()`, `summarise_at()`, `summarise_(...)`, `summarise_(..., .by_group = TRUE)`, `summarise_(..., .by_group = FALSE)`, `summarise_(..., .by_group = TRUE, .by_group_fn = function(...))`.
- Group Cases**: Use `group_by()` to "group" copy of a table. `dplyr` functions will manage each "group" separately and return the results.
- ADD CASES**: `mutate()`, `group_by()`, `ungroup()`, `group_modify()`, `group_by_(...)`, `group_modify_(...)`, `group_by_all()`, `group_modify_all()`, `group_by_if()`, `group_modify_if()`, `group_by_if_all()`, `group_modify_if_all()`.
- ADD COLUMNS**: `add_column()`, `add_new_column()`, `add_new_columns()`, `add_new_columns_(...)`, `add_new_columns_if()`, `add_new_columns_if_(...)`, `add_new_columns_if_all()`, `add_new_columns_if_all_(...)`.
- ARRANGE CASES**: `arrange()`, `arrange_(...)`, `arrange_if()`, `arrange_if_(...)`, `arrange_if_all()`, `arrange_if_all_(...)`.
- SELECT CASES**: `select()`, `select_(...)`, `select_if()`, `select_if_(...)`, `select_all()`, `select_if_all()`.
- TRANSFORM CASES**: `transmute()`, `transmute_(...)`, `transmute_if()`, `transmute_if_(...)`, `transmute_if_all()`, `transmute_if_all_(...)`.
- MAP NEW FUNCTIONS**: Three apply vectorized functions to columns. Vectorized functions take one column and return a vector of the same length as the column (see back).
- LOGIC AND BOOLEAN OPERATORS TO USE WITH FILTER()**: `<`, `<=`, `is.na()`, `%in%`, `|`, `xor()`, `>`, `>=`, `!is.na()`, `!`, `&`.

**EXTRACT CASES**

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

|  |                                                                                                                                                                                     |
|--|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|  | <code>filter(data, ...)</code> Extract rows that meet logical criteria. <code>filter(iris, Sepal.Length &gt; 7)</code>                                                              |
|  | <code>distinct(data, ..., keep = all = FALSE)</code> Remove rows with duplicate values. <code>distinct(iris, Species)</code>                                                        |
|  | <code>sample_frac(tbl, size = 1, replace = FALSE, weight = NULL, env = parent.frame())</code> Randomly sample fraction of rows. <code>sample_frac(iris, 0.5, replace = TRUE)</code> |
|  | <code>sample_n(tbl, size, replace = FALSE, weight = NULL, env = parent.frame())</code> Randomly select size rows. <code>sample_n(iris, 10, replace = TRUE)</code>                   |
|  | <code>slice(data, ...)</code> Select rows by position. <code>slice(iris, 10:15)</code>                                                                                              |
|  | <code>top_n(x, n, wt)</code> Select and order top n entries (by group if grouped data). <code>top_n(iris, 5, Sepal.Width)</code>                                                    |
|  | <code>sample_n(tbl, size, replace = FALSE, weight = NULL, env = parent.frame())</code> Randomly select size rows. <code>sample_n(iris, 10, replace = TRUE)</code>                   |
|  | <code>slice(data, ...)</code> Select rows by position. <code>slice(iris, 10:15)</code>                                                                                              |

**Logical and boolean operators to use with filter()**

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

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

33



33

# 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"**

34

17

# arrange()

35

## arrange()

Order rows by values in a column

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



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

- = Number of rows
- = Number of Columns



36

36

18

# arrange()

Order rows by values in a column

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

data frame to transform

name(s) of columns to arrange by



37

37

# 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     |



38

38

19

## 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 |



39

## Your Turn 4

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

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

40

## Pop Quiz

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

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

41

%>%

42

21

# Data Analysis Steps

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

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



43

43

# Data Analysis Steps

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



44

44

## 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)
```



45

45

## Data Analysis Steps

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



46

46

## Data Analysis Steps

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



47

## Shortcut to type %>%

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

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



48

# Scene

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

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

You decide to investigate... WITH DATA



49

49

## Your Turn 5

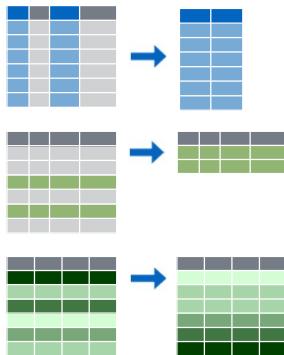
Use %>% to write a sequence of three functions that:

1. Filters to tests from the clinic (**clinic\_name**) of "picu"
2. Selects the column with the receive to verify turnaround time (**rec\_ver\_tat**) as well as the day from start of the pandemic (**pan\_day**)
3. Arrange the ` **pan\_day** ` from highest to lowest

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

50

# Isolating data



Extract variables with `select()`

Extract rows with `filter()`

Arrange rows, with `arrange()`.



51

51

# Deriving Data

52

**What is the mean and median collect to verify turnaround time by clinic?**

53

53

**Breaking down the analytical question**

1. Total TAT for each test
2. Group tests by clinic
3. Calculate mean and median for each clinic

54

54

# Deriving data



Make new variables with **mutate()**



Make summaries of data with  
**summarize()**



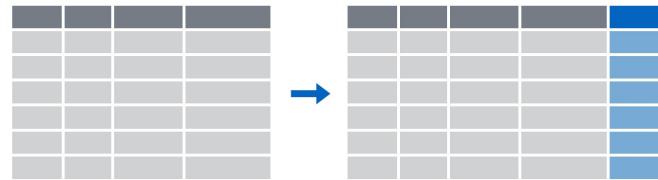
55

# mutate()

56

# mutate()

Creating new calculated columns



= Number of rows  
↑ Number of Columns



57

57

# mutate()

Creating new calculated columns

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

name for new column

equals

function whose results will populate columns



58

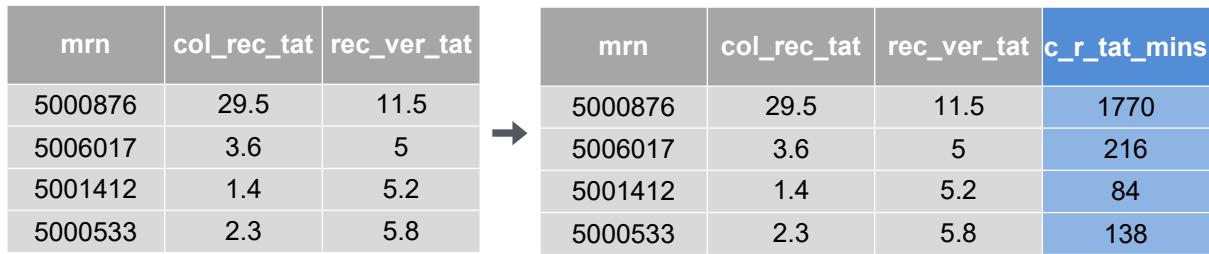
58

29

## mutate()

Creating new calculated columns

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



| 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          |

59

59

## Your Turn 6

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

60

30

# 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()  
**dplyr::cummax()** - Cumulative max()  
**dplyr::cummean()** - Cumulative mean()  
**dplyr::cummin()** - Cumulative min()  
**dplyr::cumprod()** - Cumulative prod()  
**dplyr::cumsum()** - Cumulative sum()

**RANKINGS**

**dplyr::cum\_dist()** - Proportion of all values <=  
**dplyr::dense\_rank()** - rank with ties = min, no gaps  
**dplyr::min\_rank()** - rank with ties = min  
**dplyr::nrank()** - bin into n bins  
**dplyr::percent\_rank()** - min\_rank scaled to [0,1]  
**dplyr::row\_number()** - rank with ties = "first"

**MATH**

**\*, +, ^, /, %, %%** - arithmetic ops  
**dplyr::ifelse()**, **dplyr::log10()** - 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()** - coalesce values by element across a set of vectors  
**dplyr::if\_else()** - element-wise if() + else()  
**dplyr::na\_if()** - replace specific values with NA  
**dplyr::pmax()** - element-wise max()  
**dplyr::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 groups of rows. It can also return multiple values as output.

**vectorized functions**

**COUNTS**

**dplyr::n()** - number of values/rows  
**dplyr::n\_max()** - number of unique values  
**dplyr::n\_distinct()** - number of distinct values

**LOCATION**

**dplyr::mean()**, **dplyr::median()**

**LOGICALS**

**dplyr::all()** - proportion of TRUE's  
**dplyr::is()** - check for class

**PATTERNS**

**dplyr::first()** - first value  
**dplyr::last()** - last value  
**dplyr::loc()** - location of element

**RANGE**

**dplyr::min()** - min example  
**dplyr::max()** - maximum value  
**dplyr::range()** - range example  
**dplyr::sd()** - standard deviation  
**dplyr::var()** - variance

**SPREAD**

**dplyr::gather()** - user-friendly Range  
**dplyr::pivot\_wider()** - user-friendly Range  
**dplyr::pivot\_longer()** - user-friendly Range

**WEIGHT**

**dplyr::group\_by()** - create group objects  
**dplyr::ungroup()** - remove group objects  
**dplyr::pull()** - pull values from a list  
**dplyr::select()** - select columns  
**dplyr::slice()** - slice rows  
**dplyr::slice\_head()** - slice with head/tail

**ROW NAMES**

Only data does not have rownames, which store a column of integer indices. If you want to add rownames to your data, first move them into a column.

**dplyr::rownames()**  
**dplyr::add\_rownames()**  
**dplyr::remove\_rownames()**  
**dplyr::set\_rownames()**

**Combine Tables**

**COMBINE VARIABLES**

**dplyr::bind\_rows()** - Return tables placed side by side  
**dplyr::bind\_cols()** - Return tables placed end-to-end  
BE SURE THAT ROWS ALIGN!

**COMBINE CASES**

**dplyr::bind\_rows()** - Return tables placed side by side as they are  
**dplyr::bind\_rows(..., .id = NULL)** - Return tables placed side by side as a single table. Set .id to a column name to keep the original row IDs (not preferred).

**dplyr::inner\_join()**, **dplyr::left\_join()**, **dplyr::right\_join()**, **dplyr::full\_join()** - Relatively matches

**dplyr::full\_join(..., by = NULL)** - Join matching values from x to y  
**dplyr::left\_join(..., by = NULL)** - Join matching values from x to y  
**dplyr::right\_join(..., by = NULL)** - Join matching values from y to x  
**dplyr::inner\_join(..., by = NULL)** - Relatively matches

**Extract Rows**

**dplyr::filter()** - filter one or more columns  
**dplyr::pull()** - filter one or more columns  
**dplyr::select()** - filter one or more columns  
**dplyr::slice()** - filter one or more rows  
**dplyr::subset()** - filter one or more rows  
**dplyr::unite()** - filter one or more columns  
**dplyr::unstack()** - filter one or more columns  
**dplyr::where()** - filter one or more columns

**61**

61

## Goal

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

## Objectives

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

62

62

# Dplyr tips and tricks

63

## select()

### Renaming columns

```
covid_testing %>%
  select(MRN = mrn, first_name, last_name)
```

| mrn<br><dbl> | first_name<br><chr> | last_name<br><chr> |
|--------------|---------------------|--------------------|
| 5001412      | jhezane             | westerling         |
| 5000533      | penny               | targaryen          |
| 5009134      | grunt               | rivers             |
| 5008518      | melisandre          | swyft              |

| MRN<br><dbl> | first_name<br><chr> | last_name<br><chr> |
|--------------|---------------------|--------------------|
| 5001412      | jhezane             | westerling         |
| 5000533      | penny               | targaryen          |
| 5009134      | grunt               | rivers             |
| 5008518      | melisandre          | swyft              |
|              | ..                  | ..                 |



64

64

## filter()

Filter to multiple matches

```
covid_testing %>%  
  filter(first_name %in% c("jon", "daenerys"))
```



| mrn     | first_name |
|---------|------------|
| <dbl>   | <chr>      |
| 5001412 | jhezane    |
| 5000533 | penny      |
| 5009134 | grunt      |
| 5008518 | melisandre |
| 5008967 | rolley     |

| mrn     | first_name |
|---------|------------|
| <dbl>   | <chr>      |
| 5002427 | daenerys   |
| 5011120 | jon        |
| 5001092 | jon        |
| 5004082 | jon        |
| 5005197 | daenerys   |



65

65

## mutate()

Replacing columns

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

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



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

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



66

66

# mutate()

Conditionally replacing values

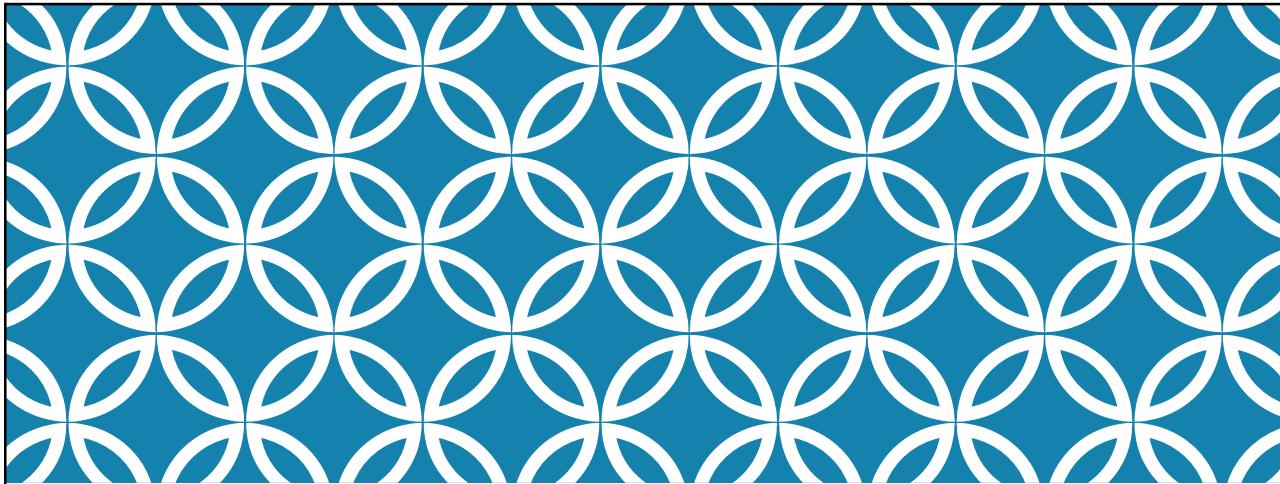
```
covid_testing %>%
  mutate(last_name = if_else(last_name=="targaryen",
                             "TARGARYEN",last_name))
```

| mrn     | first_name | last_name  |
|---------|------------|------------|
|         | <dbl>      | <chr>      |
| 5001412 | jhezane    | westerling |
| 5000533 | penny      | targaryen  |
| 5009134 | grunt      | rivers     |
| 5008518 | melisandre | swyft      |



| mrn     | first_name | last_name  |
|---------|------------|------------|
|         | <dbl>      | <chr>      |
| 5001412 | jhezane    | westerling |
| 5000533 | penny      | TARGARYEN  |
| 5009134 | grunt      | rivers     |
| 5008518 | melisandre | swyft      |





# Data Understanding: Grouping and Summarizing Data

Patrick Mathias  
May 9, 2022

1

## Goals

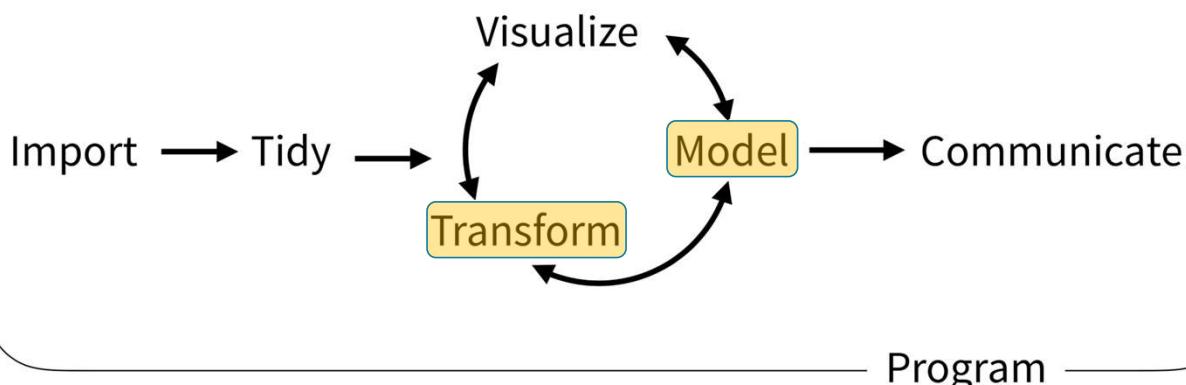
1. Learn dplyr tools for grouping and summarizing data in R

## Objectives

1. Calculate a summary statistic for a variable using the `summarize()` function
2. Creates groupings of data using the `group_by()` function
3. Combine `group_by()` and `summarize()` functions to calculate summary statistics for groups of data

2

## Typical Data Science Pipeline



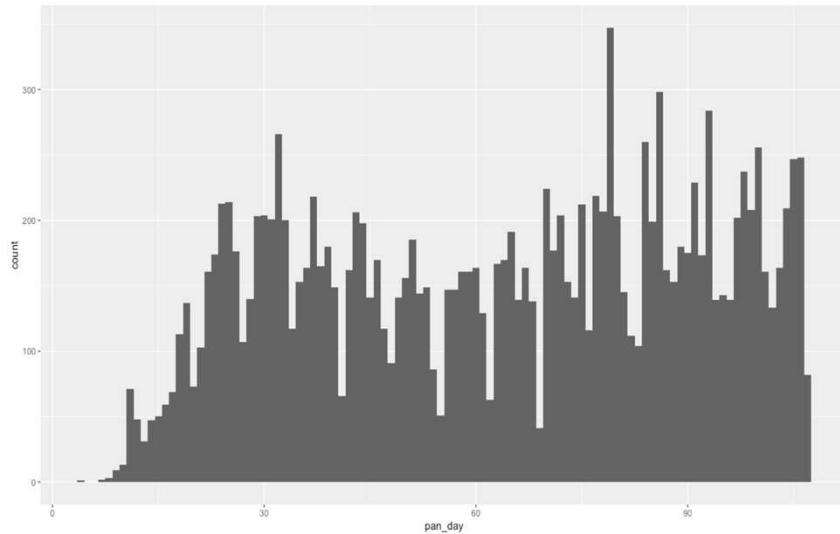
From *R for Data Science* (<https://r4ds.had.co.nz/introduction.html>)

3

Summarize the data set

4

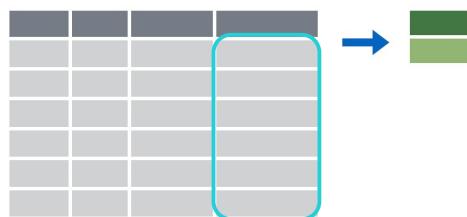
Q: How many tests are ordered per day?



5

## summarize()

- Make summaries of your data



6

## summarize()

- Make summaries of your data

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

name for new variable

Value or function

Performs calculation across all rows of data frame



7

## summarize()

- Make summaries of your data

function that returns number of observations

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

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



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



8

## summarize()

- Additional summaries = new columns

```
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         |



9

## summarize()

- Summarize supports calculations on summary stats

```
covid_testing %>%
  summarize(order_count = n(),
            day_count = n_distinct(pan_day),
            orders_per_day = order_count/day_count)
```

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



| order_count | day_count | orders_per_day |
|-------------|-----------|----------------|
| 15524       | 102       | 152            |



10

# Your Turn #1

- Open “05 – Group and Summarize.Rmd”
- Run the setup chunk
- Fill-in the gaps to calculate the mean count of orders per clinic

05:00

11

**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

**Row Names**

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first convert to a tibble.

- rownames\_to\_column()
- Move row names into col
- $\alpha < \text{rownames\_to\_column}(iris, var = "C")$
- column\_to\_rownames()
- Move col in row names.
- column\_to\_rownames(g, var = "C")

Also has\_rownames(), remove\_rownames()

**Summary Functions**

**TO USE WITH SUMMARISE()**

summary() 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

**Combine Tables**

**COMBINE VARIABLES**

Use `bind_col(x, y)` to paste tables beside each other as they are.

**COMBINE CASES**

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

**dplyr**

**COMBINE CASES**

Use `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 original table names (as pictured).

**INTERSECT**

Rows that appear in both `x` and `y`.

**SETDIF**

Rows that appear in `x` but not `y`.

**UNION**

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

**SETEQUAL**

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

**EXTRACT ROWS**

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

**SEMISEMISE**

Return rows of `x` that have a match in `y`. USEFUL TO SEE WHAT WILL BE JOINED.

**ANTI JOIN**

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

12

## Output the last day

```
covid_testing %>%
  summarize(last_day = last(pan_day))
```

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

13

## Calculate the mean turnaround time

```
covid_testing %>%
  mutate(col_ver_tat = col_rec_tat + rec_ver_tat) %>%
  summarize(col_ver_tat_mean = mean(col_ver_tat))
```

| mrn     | pan_day | col_ver_tat | col_ver_tat_mean |
|---------|---------|-------------|------------------|
| 5001412 | 4       | 6           |                  |
| 5000533 | 7       | 8           |                  |
| 5009134 | 7       | 10          |                  |
| 5008518 | 8       | 11          | 8.75             |

14

14

## Calculate the 75<sup>th</sup> percentile turnaround time

```
covid_testing %>%
  mutate(col_ver_tat = col_rec_tat + rec_ver_tat) %>%
  summarize(col_ver_tat_mean = mean(col_ver_tat),
            col_ver_75_pctile = quantile(col_ver_tat, 0.75))
```

| mrn     | pan_day | col_ver_tat | col_ver_tat_mean | col_ver_tat_mean |
|---------|---------|-------------|------------------|------------------|
| 5001412 | 4       | 6           |                  |                  |
| 5000533 | 7       | 8           |                  |                  |
| 5009134 | 7       | 10          |                  |                  |
| 5008518 | 8       | 11          | 8.75             | 9.6              |

15

15

## Your Turn #2

For the covid\_testing data frame, calculate both the median and the 95th percentile collect-to-verify turnaround time.



16

16

## Pop Quiz

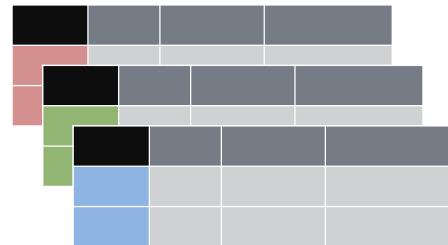
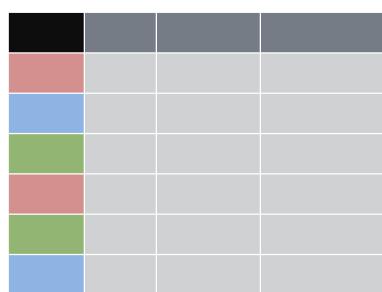
How would you calculate the median number of orders per day?

17

**Grouping your data**

18

## group\_by()



19

## group\_by()

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

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

name of variable  
to group by



20

## 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
```



21

## group\_by()

- *Group observations by `pan\_day` and `clinic\_name`*

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

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



22

`group_by() %>% summarize()`

23

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

Make summaries of your data *by group*



24

## group\_by() %>% summarize()

- Make summaries of your data

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

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



order\_count

15524



25

## 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 |



26

## Your Turn #3

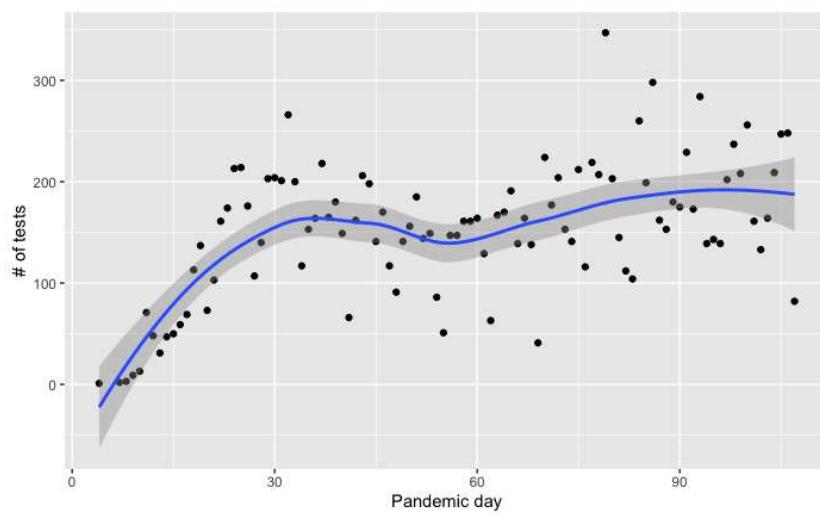
Calculate:

- a) The median collect-to-verify turnaround time for each day
- b) The median collect-to-verify turnaround time for each clinic/unit
- c) The median number of orders per day

05 : 00

27

`group_by() %>% summarize(): Example`



28

## Recap

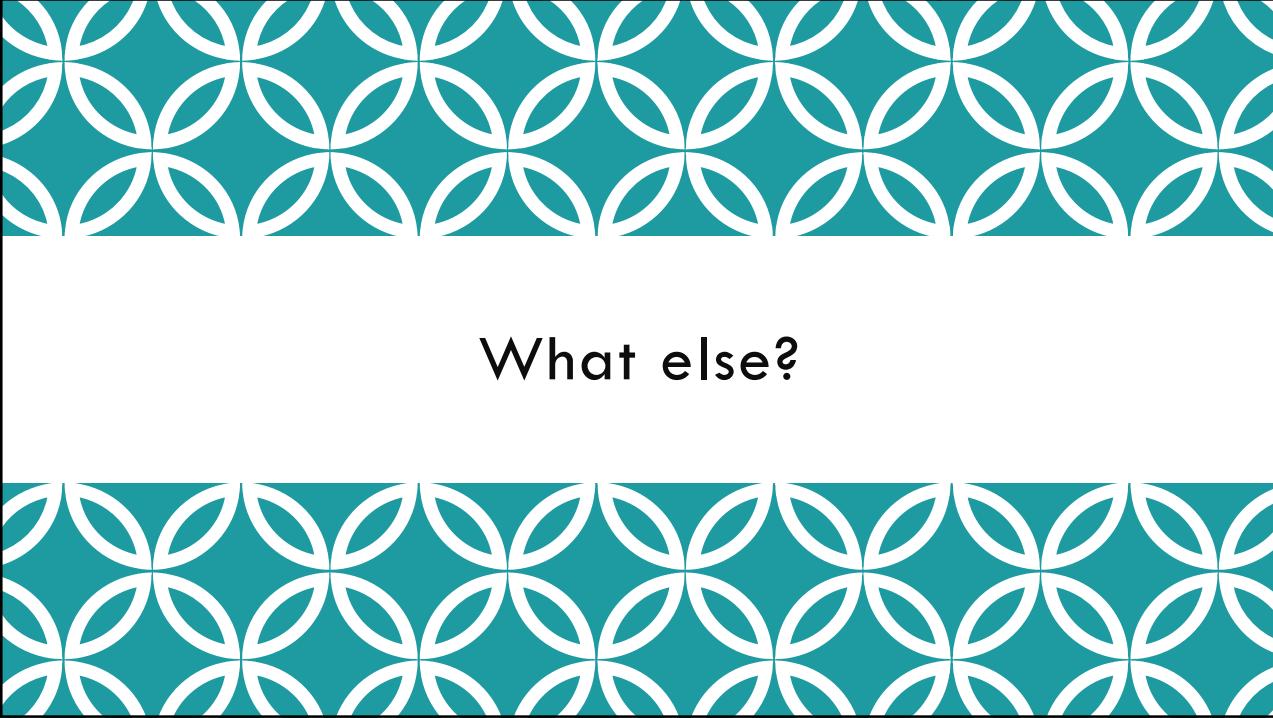
**Summarize()** is a function that enables us to calculate summaries of variables (columns).

Common summary activities include counting observations using **n()**, counting unique observations using **n\_distinct()**, and calculating means using **mean()**.

**Group\_by()** is a function that enables us to create subsets of data by a variable. Data can also be grouped by multiple variables.

Combining the **group\_by()** and **summarize()** functions is a powerful way to look at summarizations across groups.

29



What else?

30

**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

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
  → Compute table of summaries.
  summarise(mtcars, avg = mean(mpg))
```

### Group Cases

Use `group_by`(`data, ...`, `add = FALSE`, `drop = TRUE`) to create a "grouped" copy of a table grouped by columns in ... dplyr functions will manipulate each "group" separately and combine the results.

```
grouped (data, ...) → mtcars %>% group_by(cyl) %>% summarise(avg = mean(mpg))
```

### Rowwise Cases

Use `rowwise`(`data, ...`) to group data into individual rows. dplyr functions will compute results for each row. Also apply functions to list-columns. See dplyr cheat sheet for list-column workflow.

```
rowwise (data, ...) → starwars %>% rowwise() %>% mutate(lm_count = length(lfms))
```

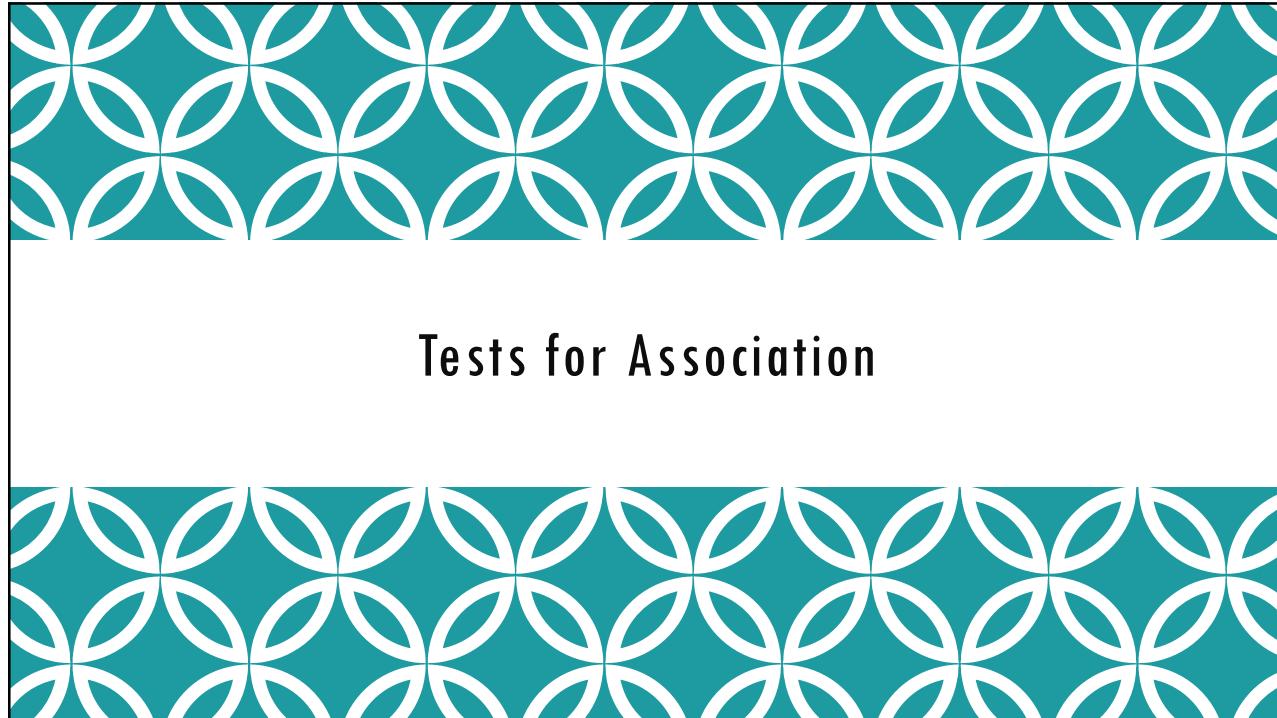
### ungroup(x, ...)

Returns ungrouped copy of table.

R Studio

RStudio® is a trademark of RStudio, PBC • CC BY SA RStudio • info@rstudio.com • 844-448-1222 • rstudio.com • Learn more at [dplyr.tidyverse.org](https://dplyr.tidyverse.org) • dplyr 1.0.7 • Updated: 2021-07

31



32

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

`payor_group_fac`  
 <chr>  
 commercial  
 government  
 other  
 unassigned

4 rows

`negative`  
 <int>  
 3549  
 3318  
 309  
 7182

`positive`  
 <int>  
 86  
 242  
 17  
 520



```
data %>%  

  fisher.test(simulate.p.value = T)
```

33

## 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"))
```



34

## 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



35

## Testing for association

| payor_group_fac |  |
|-----------------|--|
| commercial      |  |
| government      |  |
| other           |  |
| unassigned      |  |

4 rows

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

```
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
```



36

## Regression Modeling

37

Q: Is the association between test positivity and a government insurance product explained by the age of the patient?

```
tmp <- covid_testing_2 %>%
  filter(payor_group_fac %in% c("commercial", "government")) %>%
  mutate(result_fac = factor(result,
    levels=c("negative", "positive"),
    ordered=T),
    payor_group_fac = (payor_group == "government"))
tmp_fit <- glm(result_fac ~ payor_group_fac + age,      # model formula
               data = tmp,                                # dataset
               family = "binomial"                         # type of model
)
summary(tmp_fit)
exp(coefficients(tmp_fit))                                # odds
```



38

# Output for logistic regression

```

Call:
glm(formula = result_fac ~ payor_group_fac + age, family = "binomial",
     data = tmp)

Deviance Residuals:
    Min      1Q  Median      3Q     Max 
-1.6365 -0.3468 -0.2532 -0.1985  2.8393 

Coefficients:
            Estimate Std. Error z value Pr(>|z|)    
(Intercept) -4.016195  0.119611 -33.577 < 2e-16 ***
payor_group_facTRUE 1.136566  0.128761  8.827 < 2e-16 ***
age          0.032897  0.004436  7.416 1.21e-13 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 2666.6 on 7194 degrees of freedom
Residual deviance: 2535.2 on 7192 degrees of freedom
AIC: 2541.2

Number of Fisher Scoring iterations: 6

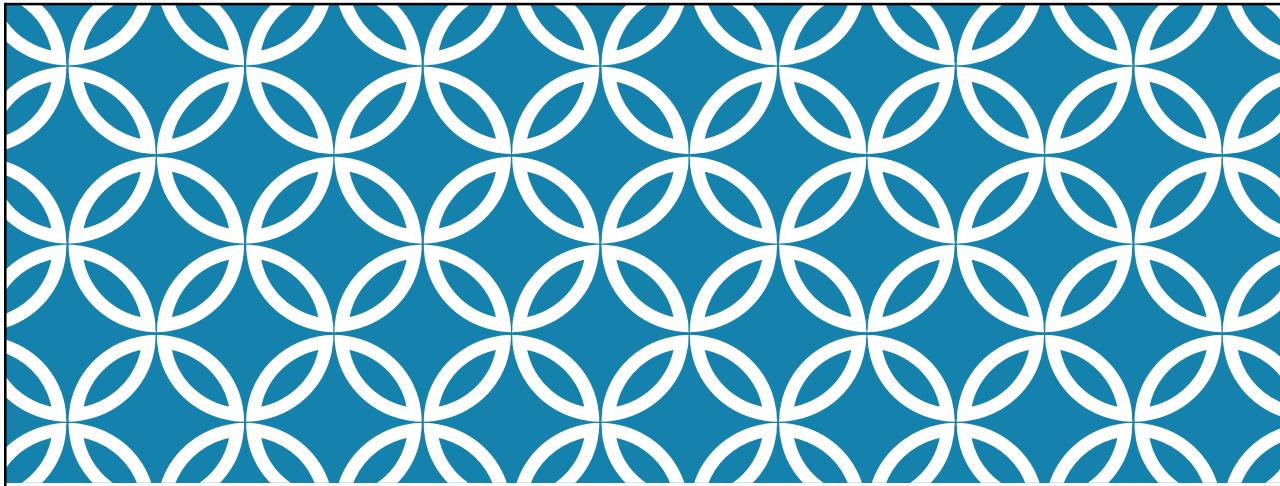
```

Odds for Payor Group  
and Age

(Intercept) payor\_group\_facTRUE  
0.01802141 3.11604971  
age  
1.03344368

39

39



# Data Understanding: Advanced Reporting

Patrick Mathias  
May 9, 2022

1

## Goals

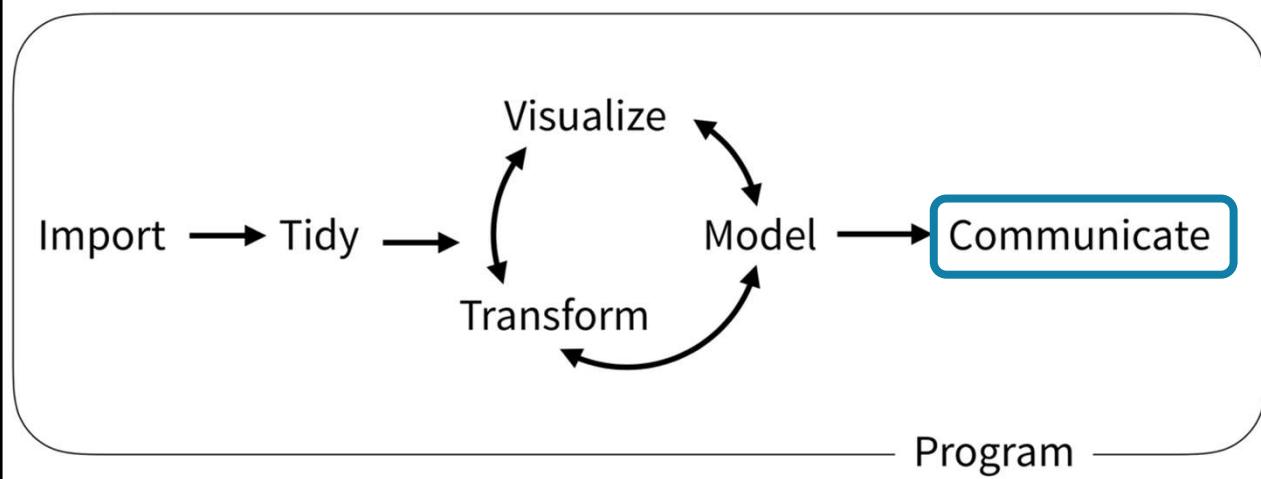
1. Build R Markdown reports using formatting outputs beyond standard document formats

## Objectives

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

2

## Typical Data Science Pipeline



3

## Refresher Quiz

You need to install 3 new packages you've never used before. What function do you run if you need to install the *flexdashboard*, *plotly*, and *DT* packages?

4

## Refresher: Installing packages

- Installing a package

```
install.packages(c("flexdashboard", "plotly", "DT"))
```

- Loading into your environment

```
library(flexdashboard)  
library(plotly)  
library(DT)
```

The flexdashboard and plotly packages were already installed in your Rstudio Cloud instance. To install them locally use `install.packages`.

5

From R Markdown to Quick and Painless Dashboards

6

**rbokeh iris dataset**

**Dashboards to drive improvement**

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

Intended to represent multiple pieces of information at a glance

7

**flexdashboard provides easy dashboard templates for reporting**

Produces HTML file that can be opened on web browsers

Or deployed on existing web server

Provides row or column based layouts

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

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

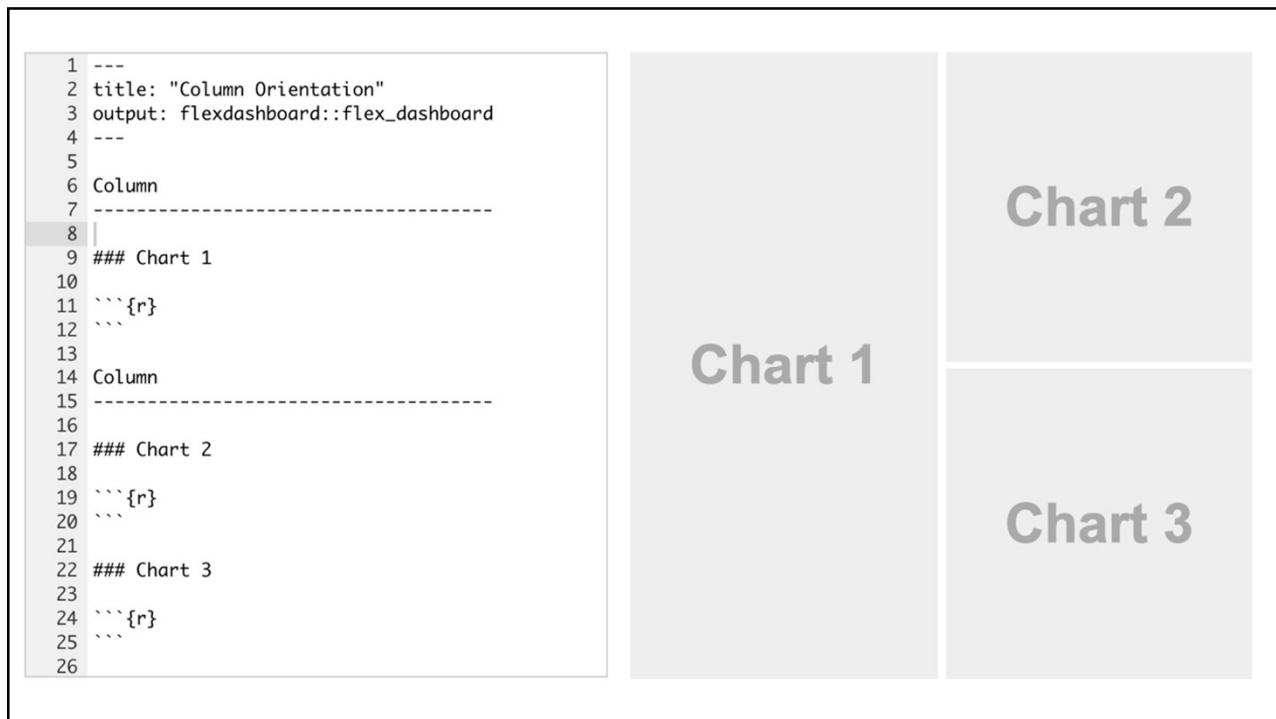
8

```

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

```

9



10

```

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

```

**Chart 1**

**Chart 2**

**Chart 3**

11

```

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

```

**Chart 1**

**Chart 2**

**Chart 3**

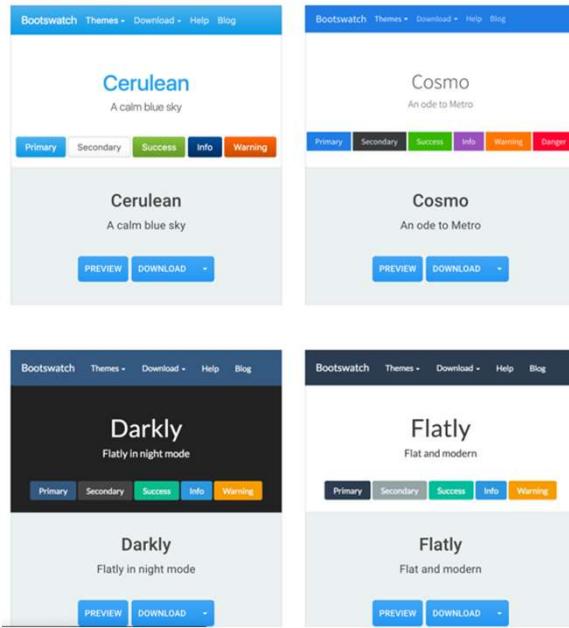
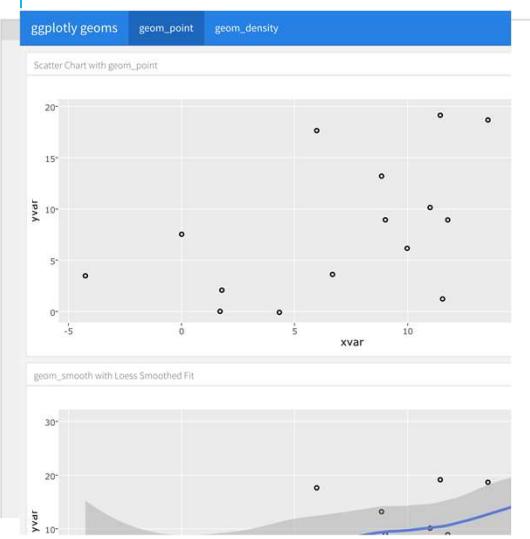
12

# Your Turn #1

1. Open “06 - Advanced Reporting.Rmd” to work with a draft COVID-19 flexdashboard and run the setup chunk. Knit the document to see the dashboard output.
2. The “Test Volumes Over Time” plot could show additional information regarding positive tests. Add fill to your barplot to show the result field in addition to overall test volume by day. Run that code chunk to see the output.
3. Too much information is crunched on the right side. Change the layout from columns to a row orientation. The 2<sup>nd</sup> and 3<sup>rd</sup> plots (Turnaround Times and Cycle Thresholds) should appear on the 2<sup>nd</sup> row.

13

## Customization



14

## Making plots interactive

15

### htmlwidgets for R support interactive visuals

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

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

16

## Plotly package converts ggplot with a simple command

To use Plotly on your desktop install the plotly package using the following command:

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

Examples of visualizations at Plotly website:

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

17

## Store plot as object and add one line to make interactive

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

18

## Your Turn #2

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

19

## Other options for interactive plots

Other interactive plot packages:

- `rbokeh`
- `Highcharter`

Time series graphs with `dygraphs` package

Maps with `leaflet` package

20

## Interactive tables with one line

DataTables library quickly converts tables into interactive element

DT package in R – to install use:

```
install.packages("DT")
```

Use datatable() function on a data frame to allow:

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

21

## datatable example

| datatable(head(iris), class = 'cell-border stripe') |  |                              |             |              |             |         |
|-----------------------------------------------------|--|------------------------------|-------------|--------------|-------------|---------|
| Show 10 ▾ entries                                   |  | Search: <input type="text"/> |             |              |             |         |
|                                                     |  | Sepal.Length                 | Sepal.Width | Petal.Length | Petal.Width | Species |
| 1                                                   |  | 5.1                          | 3.5         | 1.4          | 0.2         | setosa  |
| 2                                                   |  | 4.9                          | 3           | 1.4          | 0.2         | setosa  |
| 3                                                   |  | 4.7                          | 3.2         | 1.3          | 0.2         | setosa  |
| 4                                                   |  | 4.6                          | 3.1         | 1.5          | 0.2         | setosa  |
| 5                                                   |  | 5                            | 3.6         | 1.4          | 0.2         | setosa  |
| 6                                                   |  | 5.4                          | 3.9         | 1.7          | 0.4         | setosa  |

Showing 1 to 6 of 6 entries      Previous 1 Next

22

## Your Turn #3

We are going to replace one of our panels of content with an interactive table.

1. Load the DataTables package in your setup chunk
2. Rename the cycle threshold distribution panel to “Positive Result Details”. Use filter() to create a dataframe that only includes positive results.
3. Display an interactive table that includes the content from the positive result dataframe you created.
4. Knit the dashboard. Search “lannister” in the search box to confirm the interactive table is working.

23

## Your Turn #4

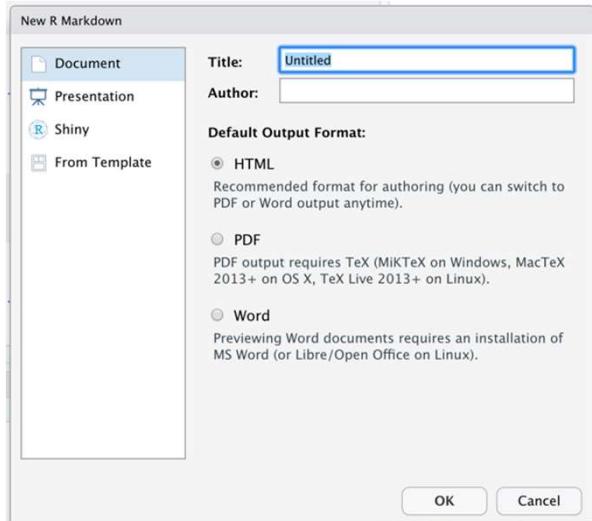
Customize your dashboard. Use any of the data available in your covid testing dataset to generate new plots or tables that provide insight into the underlying data.

24

## What Else?

25

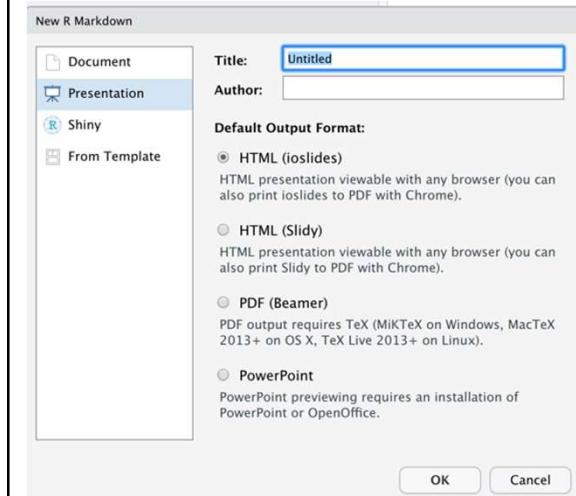
### Standard Markdown Reporting Formats



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

26

## Formats to go straight from code to slides



Multiple HTML formats create webpage that's advanceable like slides

PDF presentation uses LaTeX in the background

Powerpoint produces simple slides

27

```

1 ---  

2 title: "Untitled"  

3 output: powerpoint-presentation ← Output format  

4 ---  

5  

6 ```{r setup, include=FALSE}  

7 knitr::opts_chunk$set(echo = FALSE)  

8 ```  

9  

10 ## R Markdown  

11  

12 This is an R Markdown presentation. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see http://rmarkdown.rstudio.com.  

13  

14 When you click the Knit button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document.  

15  

16 ## Slide with Bullets ← Each slide has its own header  

17  

18 - Bullet 1  

19 - Bullet 2  

20 - Bullet 3  

21  

22 ## Slide with R Output  

23  

24 ```{r cars, echo = TRUE}  

25 summary(cars)  

26 ```  

27  

28 ## Slide with Plot  

29  

30 ```{r pressure}  

31 plot(pressure)  

32 ```  

33  

34

```

28

## R Markdown

This is an R Markdown presentation. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document.

Example output slide

29

## Slide with Bullets

- Bullet 1
- Bullet 2
- Bullet 3

Example output slide

30

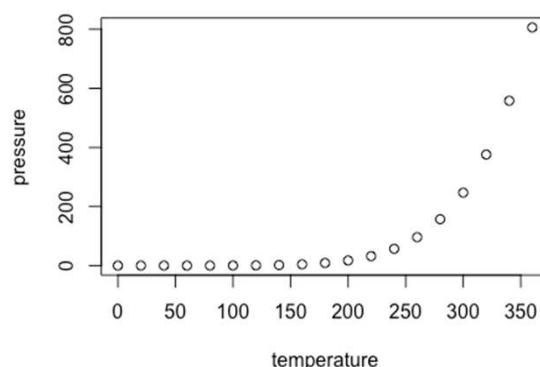
## Slide with R Output

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

Example output slide

31

## Slide with Plot



Example output slide

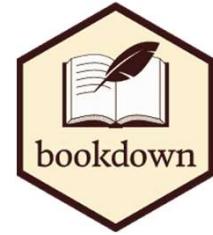
32

## Books and longer documents also generated from R Markdown

Can generate printer ready books and ebooks

Supports LaTeX features such as equations

Generates blog formatted websites



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

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

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

33

## Goals

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

## Objectives

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

34

# Appendix

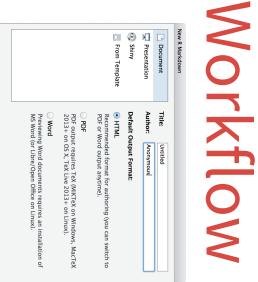
# R Markdown :: CHEAT SHEET

## 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

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

- Write document by editing template

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

- Preview Output in IDE window

- Publish (optional) to web server

- Examine build log in R Markdown console

- Use output file that is saved alongside .Rmd

## Embed code with knitr syntax

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

**Built with r getVersion()** ↘ Built with 3.2.3

getVersion()

## [1] '3.2.3'

## 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 = '##')

**eval** - Run code in chunk (default = TRUE)

Options not listed above: R.options, aniopts, autoprop, 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, purr, ref.label, render, size, split, tidy.opts



## .rmd Structure

**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`

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**

ends with `r`

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**

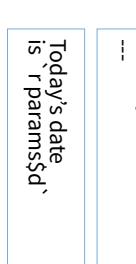
For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

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

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

3. **Set parameters** • Set values with Knit with parameters or the params argument with render():

render("doc.Rmd", params = list(n = 1, d = as.Date("2015-01-01")))



## Parameters

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

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

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

3. **Set parameters** • Set values with Knit with parameters or the params argument with 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

- Add runtime: shiny to the YAML header.

- Call Shiny input functions to embed input objects.

- 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}

if (echo = FALSE)

numericInput("n",

"How many cars?", 5)

renderTable({

head(cars, input\$n)}

})



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)



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.

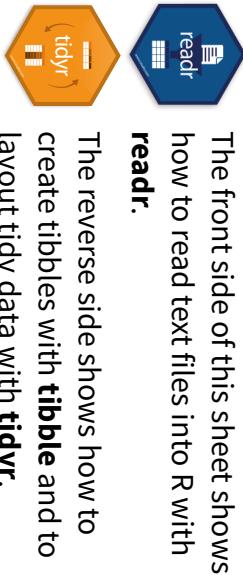




# Data Import :: CHEAT SHEET

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**.

## 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())
```

## 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.

**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
- **rvest** - HTML (Web Scraping)

- **tidy** -

create tibbles with **tibble** and to layout tidy data with **tidyr**.

### Comma Delimited Files

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

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

To make file.csv run:

write\_file(x = "a,b,c\n1,2,3\n4,5,NA", path = "file.csv")

### Semi-colon Delimited Files

```
a;b;c  
1;2;3  
4;5;NA
```

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

To make file2.csv run:

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

### Files with Any Delimiter

```
a|b|c  
1|2|3  
4|5|NA
```

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

To make file.txt run:

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

### Fixed Width Files

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

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

To make file.fwf run:

write\_file(x = "a\tb\tc\n1\t2\t3\n4\t5\tNA", path = "file.fwf")

### Tab Delimited Files

read\_tsv("file.tsv") Also **read\_table()**.

write\_file(x = "a\tb\tc\t1\t2\t3\t4\t5\tNA", path = "file.tsv")

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 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())



## Data types

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

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

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)

use **problems()** to diagnose problems.

**X <- read\_csv("file.csv"); problems(X)**

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()))**

Else, read in as character vectors then parse with a **parse\_** function.

**• parse\_guess()**

**• parse\_character()**

**• parse\_datetime()** Also **parse\_date()** and

**• 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 x 6
  manufacturer model disp    cyl  year  trans
  <chr>        <chr> <dbl> <dbl> <dbl> <chr>
  1 audi         a4      1.8   4     1999 auto
  2 audi         a4      1.8   4     1999 auto
  3 audi         a4      1.8   4     1999 auto
  4 audi         a4      1.8   4     1999 auto
  5 audi         a4      1.8   4     1999 auto
  6 audi         a4      1.8   4     1999 auto
  7 audi         a4      1.8   4     1999 auto
  8 audi         a4      1.8   4     1999 auto
  9 audi         a4      1.8   4     1999 auto
  # ... with 244 more rows, and 3 more variables:
  #   cyl <dbl>, trans <chr>, year <dbl>
  #   note: more variables: year <dbl>
  #   note: omitted 68 rows]
```

## tibble display

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

table4a

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

- 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)**
- spread(data, key, value, fill = NA, convert = FALSE, drop = TRUE, sep = NULL)**
- gather()** moves column names into a **key** column, gathering the column values into a single value column.

table2

| country | year | type  | count |
|---------|------|-------|-------|
| A       | 1999 | cases | 0.7K  |
| A       | 1999 | pop   | 19M   |
| B       | 1999 | cases | 2K    |
| B       | 1999 | pop   | 172M  |
| C       | 2000 | cases | 80K   |
| C       | 2000 | pop   | 174M  |

- 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.

table3

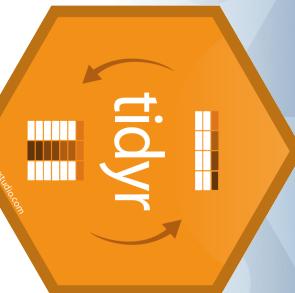
| country | year | rate     |
|---------|------|----------|
| A       | 1999 | 0.7K/19M |
| A       | 2000 | 2K/20M   |
| B       | 1999 | 37K/172M |
| B       | 2000 | 80K/174M |
| C       | 1999 | 212K/1T  |
| C       | 2000 | 213K/1T  |

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

table4

| 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  |

# 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", ...)**

table5

| country | century | year |
|---------|---------|------|
| Afghan  | 19      | 99   |
| Afghan  | 20      | 00   |
| Brazil  | 19      | 99   |
| Brazil  | 20      | 00   |
| China   | 19      | 99   |
| China   | 20      | 00   |
| China   | 1999    | 2000 |

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

table6

unite(table6, century, year, col = "year", sep = "")

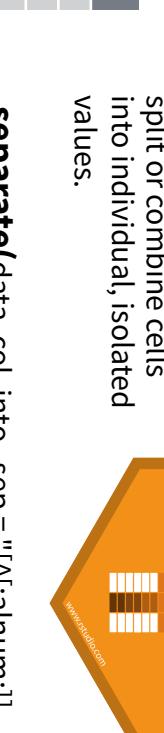
table7

unite(table7, century, year, col = "year", sep = "")

# Tidy Data with tidyverse

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

A table is tidy if:



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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

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

- Each **variable** is in its own **column**
- Each **observation**, or **case**, is in its own **row**
- Makes variables easy to access as vectors
- Preserves cases during vectorized operations

- <ul

# 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**().

**distinct**(data, ..., keep\_all = FALSE) Remove rows with duplicate values.

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

**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 =  
 ungroup(x, ...))`

`FALSE)`  
Returns copy of table grouped by ...  
`g Iris <- group_by(iris, Species)`

## 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.

**sample\_frac**(tbl, size = 1, replace = FALSE, weight = NULL, .env = parent.frame(), select\_fraction(iris, 0.5, replace = TRUE)

**sample\_n**(tbl, size, replace = FALSE, weight = NULL, .env = parent.frame(), 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)`

**arrange**(.data, ..., .by = ..., .order = "asc") Order rows by values of a column or columns (low to high), use with **desc**() to order from high to low.

**arrange**(mtcars, mpg)  
`arrange(mtcars, desc(mpg))`

**mutate**(.data, ...) Order rows by values of a column or columns (low to high), use with **desc**() to order from high to low.

**mutate**(.tbl, .funs, ...) Apply funs to every column. Use with **funs**() (Also **mutate\_if**()).  
`mutate_all(faithful, funs(log(.), log2(.)))`

**mutate**(.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 CASES**

**add\_row**(.data, ..., before = NULL, after = NULL) Add one or more rows to a table.  
`add_row(faithful, eruptions = 1, waiting = 1)`

**add\_column**(.data, ..., before = NULL, after = NULL) Add new column(s). Also **add\_count**(), **add\_tally**(). `add_column(mtcars, new = 1:32)`

### EXTRACT VARIABLES

## Manipulate Variables

### EXTRACT CASES

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\_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\_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(.)))`

**RENAME CASES**

**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()  
`dplyr::cummax()` - Cumulative max()  
`dplyr::cummean()` - Cumulative mean()  
`cummin()` - Cumulative min()  
`cumprod()` - Cumulative prod()  
`cumsum()` - Cumulative sum()

**RANKINGS**  
`dplyr::cume_dist()` - Proportion of all values <= rank wties = min, no gaps  
`dplyr::dense_rank()` - rank wties = min  
`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"

**CUMULATIVE MEANS**  
`cummean()` - Cumulative mean()

**CUMULATIVE MIN**  
`cummin()` - Cumulative min()

**CUMULATIVE PROD**  
`cumprod()` - Cumulative prod()

**CUMULATIVE SUM**  
`cumsum()` - Cumulative sum()

**LOCATION**  
`dplyr::n()` - number of values/rows  
`dplyr::n_distinct()` - # of uniques  
`sum(!is.na())` - # of non-NAs

**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

**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",`  
`TRUE ~ Species))`

**SPREAD**  
`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-NAs

**MEANS**  
`mean()` - mean, also **mean(is.na())**  
`median()` - median

**POSITION**  
`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

**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",`  
`TRUE ~ Species))`

**SPREAD**  
`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`

# Combine Tables

## COMBINE VARIABLES

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

## COMBINE CASES

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   1 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |

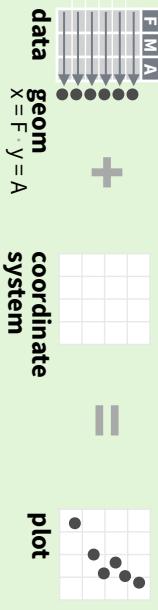
|           |   |           |
|-----------|---|-----------|
| X         |   | Y         |
| A   B   C |   | A   B   C |
| a   t   1 | + | a   t   3 |
| b   u   2 | = | b   u   2 |
| c   v   3 |   | c   v   3 |



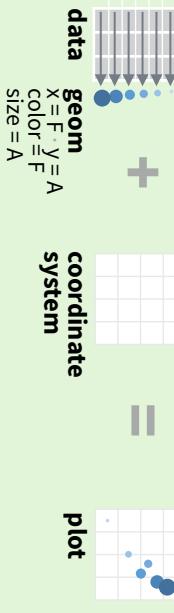
# Data Visualization with ggplot2 :: CHEAT SHEET

## 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.

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

a + geom_point()
#(Useful for expanding limits)

b + geom_curve(aes(yend = lat + 1, xend = long + 1), curvature = -1)
# 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 + 1, ymin = lat, xmax =
  long + 1, ymax = lat + 1))
# x, y, alpha, color, fill, linetype, size

a + geom_ribbon(aes(ymin = unemploy + 900,
  ymax = unemploy + 900))
# x, y, max, 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")
c + geom_density(kernel = "gaussian")
c + geom_dotplot()
x, y, alpha, color, fill, linetype, size
```

```
c + geom_hex()
x, y, alpha, color, fill, linetype, size
```

```
c + geom_rug(sides = "bl")
x, y, alpha, color, fill, linetype, size
```

```
d + geom_smooth(method = lm)
x, y, alpha, color, fill, group, linetype, size
```

```
e + geom_text(label = "text", nudge_x = 1, check_overlap = TRUE)
x, y, alpha, color, fill, shape, size, stroke
```

```
f + geom_raster(aes(fill = z), interpolate = FALSE)
x, y, alpha, fill
```

```
g + geom_bar(aes(x, alpha, color, fill, linetype, size, weight))
x, alpha, color, fill, linetype, size, weight
```

```
ggplot(data = mpg, aes(x = cty, y = hwy)) + geom_point()
# 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

**gg\_save("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))
```

```
c <- ggplot(mpg, aes(cty, hwy))
```

```
d <- ggplot(diamonds, aes(cut, color))
```

```
e <- ggplot(diamonds, aes(cut, color))
```

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

```
g <- ggplot(mpg, aes(cty, hwy))
```

```
h <- ggplot(diamonds, aes(carat, price))
```

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

```
j <- ggplot(df, aes(grp = c("A", "B"), fit = 4.5, se = 1.2))
```

```
k <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
l <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
m <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
n <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
o <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
p <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
q <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
r <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
s <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
t <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
u <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
v <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
w <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
x <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
y <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

```
z <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

### TWO VARIABLES

#### continuous x, continuous y

```
e + geom_label(aes(label = cty), nudge_x = 1, check_overlap = TRUE)
```

```
f + geom_text(aes(label = carat), nudge_y = 1, check_overlap = TRUE)
```

```
g + geom_raster(aes(fill = z), interpolate = FALSE)
```

```
h + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

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

```
j + geom_crossbar(aes(x, y, ymin, ymax, alpha, color, fill, group, linetype, size))
```

```
k + geom_errorbar(aes(x, y, ymin, ymax, alpha, color, fill, group, linetype, size))
```

```
l + geom_errorbarh(aes(x, y, ymin, ymax, alpha, color, fill, group, linetype, size))
```

```
m + geom_pointrange(aes(x, y, ymin, ymax, alpha, color, fill, group, linetype, size))
```

```
n + geom_linerange(aes(x, y, ymin, ymax, alpha, color, fill, group, linetype, size))
```

```
o + geom_raster(aes(fill = z), interpolate = FALSE)
```

```
p + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

```
q + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

```
r + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

```
s + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

```
t + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

```
u + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

```
v + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

```
w + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

```
x + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

```
y + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

```
z + geom_hex(aes(x, y, alpha, colour, group, linetype, size))
```

### continuous bivariate distribution

```
h + geom_bin2d(binwidth = c(0.25, 500))
```

```
h + geom_hex(x, y, alpha, colour, group, linetype, size)
```

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

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

### THREE VARIABLES

```
maps
```

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

```
state = tolower(rownames(USArrests)))
```

```
map <- map_data("state")
```

```
k <- ggplot(data, aes(fill = murder))
```

```
l + geom_raster(aes(fill = z), hijust = 0.5, vjust = 0.5,
```

```
interpolate = FALSE)
```

```
x, y, alpha, fill
```

```
l + geom_map(aes(map_id = state), map = map)
```

```
+ expand_limits(x = map$long, y = map$lat),
```

```
map_id, alpha, color, fill, linetype, size)
```

### ggplot2

www.rstudio.com



## 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.