

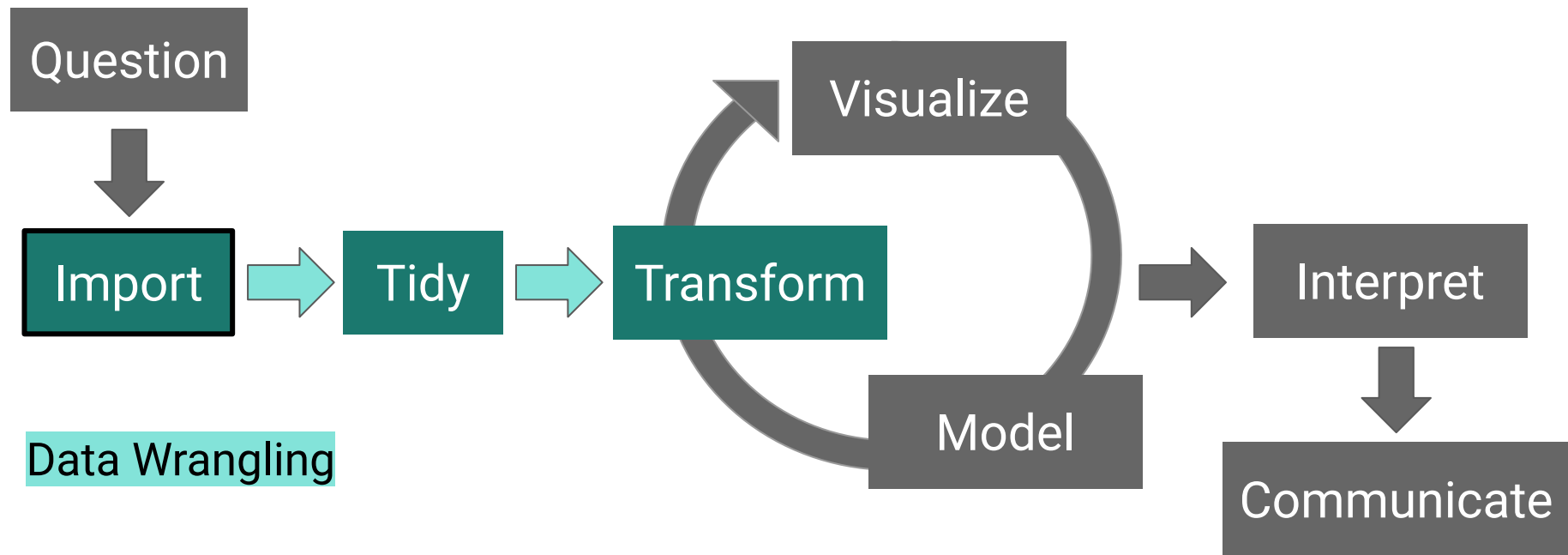
# Importing Data

## Lecture 4

# Objectives

- To import files
- To download files
- To write to a file

# Motivation



# base R vs package

- R packages are collections of functions developed by the community
- R packages improve existing base R functionalities, or by add new ones

```
# Read documentation  
> packageDescription("tidyverse")  
> help(package = "tidyverse")
```

Package: tidyverse

Title: Easily Install and Load the 'Tidyverse'

Version: 1.3.0

Authors@R: c(person(given = "Hadley", family = "Wickham", role = c("aut", "cre"), email = "hadley@rstudio.com"), person(given = "RStudio", role = c("cph", "fnd")))

Description: The 'tidyverse' is a set of packages that work in harmony because they share common data representations and 'API' design. This package is designed to make it easy to install and load multiple 'tidyverse' packages in a single step. Learn more about the 'tidyverse' at <<https://tidyverse.org>>. ...

# Tidyverse

Tidyverse

Packages

Blog

Learn

Help

Contribute



## R packages for data science

The tidyverse is an opinionated **collection of R packages** designed for data science. All packages share an underlying design philosophy, grammar, and data structures.

Install the complete tidyverse with:

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

<https://www.tidyverse.org/>

# Tibble vs Dataframe

Tidyverse

Packages



## tibble

tibble is a modern re-imagining of the data frame, keeping what time has proven to be effective, and throwing out what it has not. Tibbles are data.frames that are lazy and surly: they do less and complain more forcing you to confront problems earlier, typically leading to cleaner, more expressive code. [Go to docs...](#)

<https://www.tidyverse.org/packages/>

# readr package

Tidyverse

Packages



## readr

readr provides a fast and friendly way to read rectangular data (like csv, tsv, and fwf). It is designed to flexibly parse many types of data found in the wild, while still cleanly failing when data unexpectedly changes. [Go to docs...](#)

<https://www.tidyverse.org/packages/>

# Load tidyverse

```
> library(tidyverse)
```

— Attaching packages — tidyverse 1.3.0 —

✓ ggplot2 3.3.2 ✓ purrr 0.3.4

✓ tibble 3.0.4 ✓ dplyr 1.0.2

✓ tidyr 1.1.2 ✓ stringr 1.4.0

✓ readr 1.4.0 ✓ forcats 0.5.0

— Conflicts — tidyverse\_conflicts() —

x dplyr::filter() masks stats::filter()

x dplyr::lag() masks stats::lag()



# Importing text files

- open plain text files and convert to data frames
- text files
  - comma-separated value (csv) files
  - tab-delimited files
- other delimiters
  - semi-colon (;)
  - forward slash (/)

# Importing text files

patient.csv

patient\_ID,sex,age\_year,weight\_kg,height\_cm

P001,female,1,9.1,73

P002,female,4,16.4,96

P003,female,2,10.5,85

P004,male,3,13.2,95

P005,male,4,15.9,104

```
> patient
```

	patient_ID	sex	age_year	weight_kg	height_cm
1	P001	female	1	9.1	73
2	P002	female	4	16.4	96
3	P003	female	2	10.5	85
4	P004	male	3	13.2	95
5	P005	male	4	15.9	104

# read\_csv( )

```
> patient <- read_csv("patient.csv")
```

— Column specification —

```
cols(  
  patient_ID = col_character(),  
  sex = col_character(),  
  age_year = col_double(),  
  weight_kg = col_double(),  
  height_cm = col_double()  
)
```

```
> patient
```

# A tibble: 5 x 5

patient_ID	sex	age_year	weight_kg	height_cm
<chr>	<chr>	<dbl>	<dbl>	<dbl>
1 P001	female	1	9.1	73
2 P002	female	4	16.4	96
3 P003	female	2	10.5	85
4 P004	male	3	13.2	95
5 P005	male	4	15.9	104

# Parsing a file

- **readr** automatically guesses the type of vector in each column
- uses the **parse\_\*( )** function
  - \* = **logical, integer, double, character, date, factor**

```
> parse_double(c("123", "456", "abc", "123.456"))
```

```
[1] 123.000 456.000    NA 123.456
```

```
attr("problems")
```

```
# A tibble: 1 x 4
```

	row	col	expected	actual
		<int>	<int>	<chr>
1	3	NA	a double	abc

# read\_tsv( )

patient.txt

patient_ID	sex	age_year	weight_kg	height_cm
P001	female	1	9.1	73
P002	female	4	16.4	96
P003	female	2	10.5	85
P004	male	3	13.2	95
P005	male	4	15.9	104

```
> patient <- read_tsv("patient.txt")  
> patient
```

# A tibble: 5 x 5

	patient_ID	sex	age_year	weight_kg	height_cm
	<chr>	<chr>	<dbl>	<dbl>	<dbl>
1	P001	female	1	9.1	73
2	P002	female	4	16.4	96
3	P003	female	2	10.5	85
4	P004	male	3	13.2	95
5	P005	male	4	15.9	104

# No headers

patient\_headless.csv

P001,female,1,9.1,73  
P002,female,4,16.4,96  
P003,female,2,10.5,85  
P004,male,3,13.2,95  
P005,male,4,15.9,104

```
> patient <- read_csv("patient_headless.csv",  
  col_names = FALSE)
```

```
> patient
```

```
# A tibble: 5 x 5
```

	X1	X2	X3	X4	X5
	<chr>	<chr>	<dbl>	<dbl>	<dbl>
1	P001	female	1	9.1	73
2	P002	female	4	16.4	96
3	P003	female	2	10.5	85
4	P004	male	3	13.2	95
5	P005	male	4	15.9	104

# Provide headers

```
> headers <- c("patient", "sex", "age", "weight", "height")  
> patient <- read_csv("patient_headless.csv",  
                      col_names = headers)  
> patient
```

# A tibble: 5 x 5

	patient	sex	age	weight	height
	<chr>	<chr>	<dbl>	<dbl>	<dbl>
1	P001	female	1	9.1	73
2	P002	female	4	16.4	96
3	P003	female	2	10.5	85
4	P004	male	3	13.2	95
5	P005	male	4	15.9	104

# Skip lines

```
> patient <- read_csv("patient_headless.csv",  
  col_names = FALSE,  
  skip = 2)  
  
> patient
```

# A tibble: 3 x 5

	X1	X2	X3	X4	X5
	<chr>	<chr>	<dbl>	<dbl>	<dbl>
1	P003	female	2	10.5	85
2	P004	male	3	13.2	95
3	P005	male	4	15.9	104



# Importing Excel files

- different sheets with tabular data

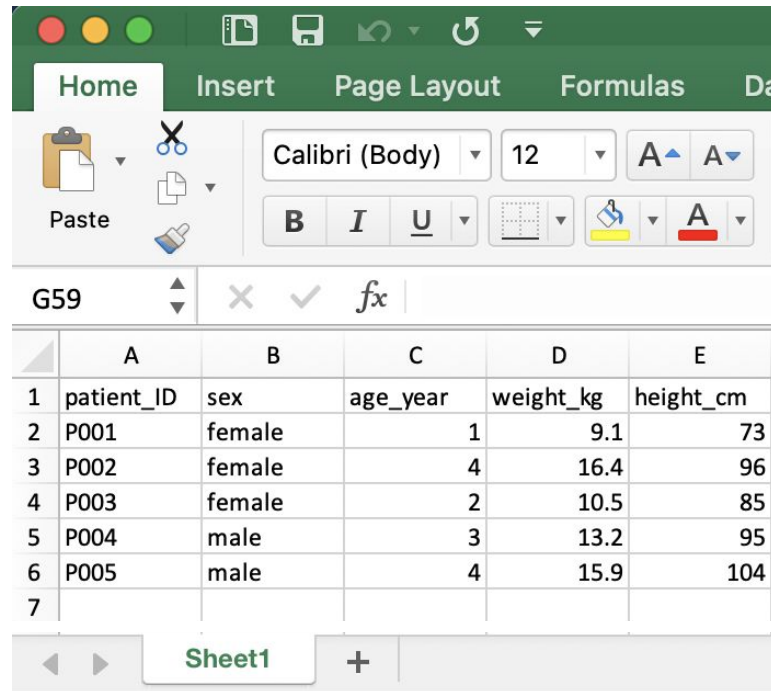
```
> install.packages("readxl")
```

```
> library(readxl)
```

```
# list different sheets
```

```
> excel_sheets("patient.xlsx")
```

```
1] "Sheet1"
```



The screenshot shows the Microsoft Excel interface. The 'Home' tab is selected in the ribbon. The font is 'Calibri (Body)' and the size is '12'. The active cell is 'G59'. The spreadsheet contains a table with 5 columns: 'patient\_ID', 'sex', 'age\_year', 'weight\_kg', and 'height\_cm'. The data rows are numbered 1 to 7. The bottom of the window shows the 'Sheet1' tab.

	A	B	C	D	E
1	patient_ID	sex	age_year	weight_kg	height_cm
2	P001	female	1	9.1	73
3	P002	female	4	16.4	96
4	P003	female	2	10.5	85
5	P004	male	3	13.2	95
6	P005	male	4	15.9	104
7					

# read\_excel( )

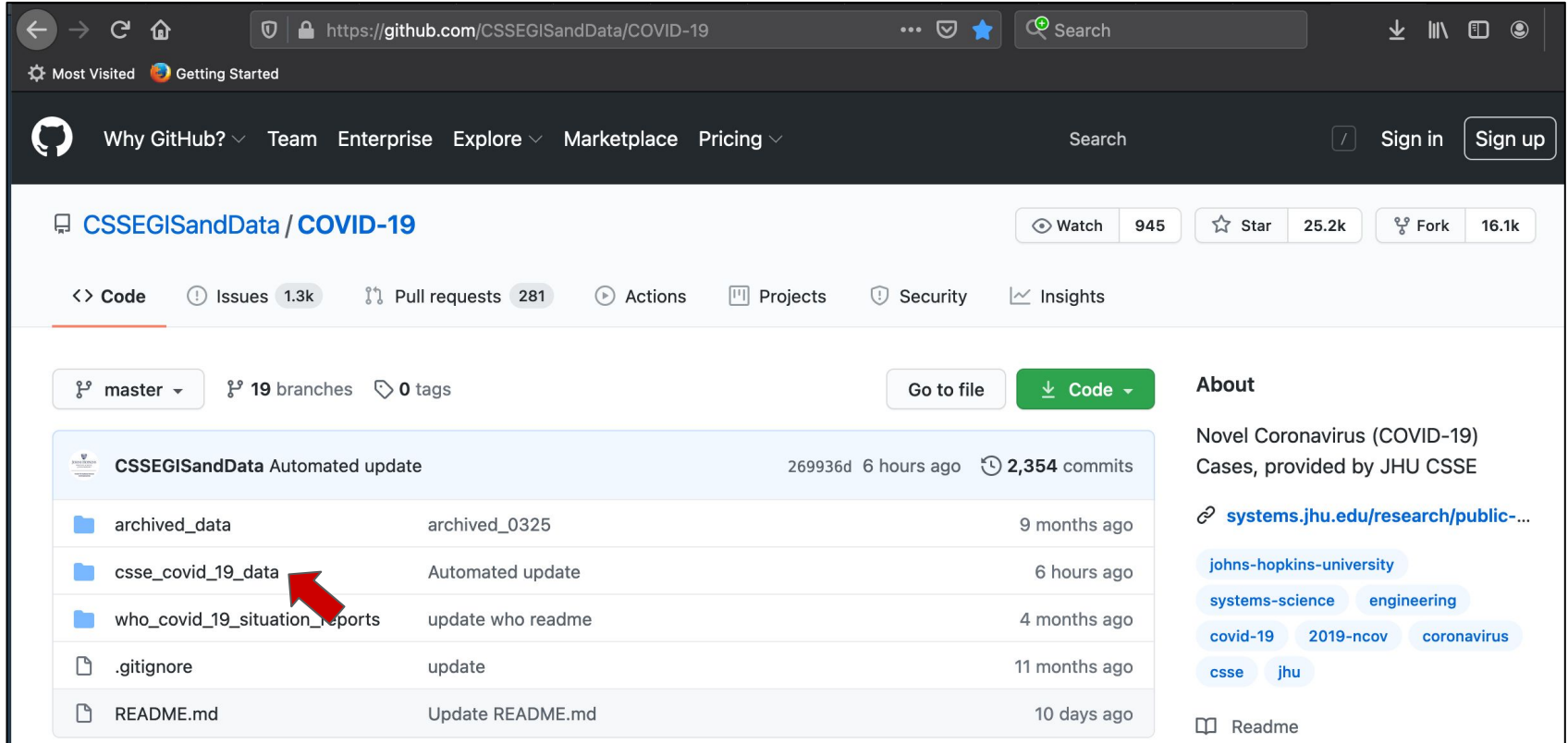
```
> patient <- read_excel("patient.xlsx", sheet = "Sheet1") # name of sheet  
> patient <- read_excel("patient.xlsx", sheet = 1)       # sheet number
```

```
> patient
```

# A tibble: 5 x 5

patient_ID	sex	age_year	weight_kg	height_cm
<chr>	<chr>	<dbl>	<dbl>	<dbl>
1 P001	female	1	9.1	73
2 P002	female	4	16.4	96
3 P003	female	2	10.5	85
4 P004	male	3	13.2	95
5 P005	male	4	15.9	104

# Importing files from the internet



The screenshot shows the GitHub repository page for **CSSEGISandData / COVID-19**. The browser address bar displays <https://github.com/CSSEGISandData/COVID-19>. The repository page includes navigation links (Why GitHub?, Team, Enterprise, Explore, Marketplace, Pricing), a search bar, and buttons for Sign in and Sign up. The repository statistics show 945 Watchers, 25.2k Stars, and 16.1k Forks. The file list includes:

File/Folder	Last Commit	Time Ago
archived_data	archived_0325	9 months ago
csse_covid_19_data	Automated update	6 hours ago
who_covid_19_situation_reports	update who readme	4 months ago
.gitignore	update	11 months ago
README.md	Update README.md	10 days ago

A red arrow points to the **csse\_covid\_19\_data** folder. The right sidebar contains an "About" section with the text: "Novel Coronavirus (COVID-19) Cases, provided by JHU CSSE" and a link to [systems.jhu.edu/research/public...](https://systems.jhu.edu/research/public...). Below this are tags for **johns-hopkins-university**, **systems-science**, **engineering**, **covid-19**, **2019-ncov**, **coronavirus**, **csse**, and **jhu**.

# Importing files from the internet

The screenshot shows the GitHub interface for the repository **CSSEGISandData / COVID-19**. The browser address bar displays the URL `https://github.com/CSSEGISandData/COVID-19/tree/master/csse_covid_19`. The repository page includes navigation links (Code, Issues, Pull requests, Actions, Projects, Security, Insights) and a file tree view. The file tree shows the following files and folders:

- `csse_covid_19_daily_reports` (Automated update, 9 hours ago)
- `csse_covid_19_daily_reports_us` (Automated update, 9 hours ago)
- `csse_covid_19_time_series` (Automated update, 9 hours ago)
- `README.md` (Update README.md, 6 days ago)
- `UID_ISO_FIPS_LookUp_Table.csv` (add Samoa, 28 days ago)

A red arrow points to the `csse_covid_19_time_series` file. The commit history for this file shows a commit by **CSSEGISandData** with the message "Automated update" at `e031127` 9 hours ago.

# Importing files from the internet

The screenshot shows a web browser displaying the GitHub repository for CSSEGISandData/COVID-19. The browser's address bar shows the URL [https://github.com/CSSEGISandData/COVID-19/tree/master/csse\\_covid\\_19\\_time\\_series](https://github.com/CSSEGISandData/COVID-19/tree/master/csse_covid_19_time_series). The repository page includes a header with the repository name, a navigation bar with links to Code, Issues, Pull requests, Actions, Projects, Security, and Insights, and a sidebar with repository statistics (Watch, Star, Fork). The main content area shows the file list for the `csse_covid_19_time_series` directory. A red arrow points to the file `time_series_covid19_confirmed_global.csv`.

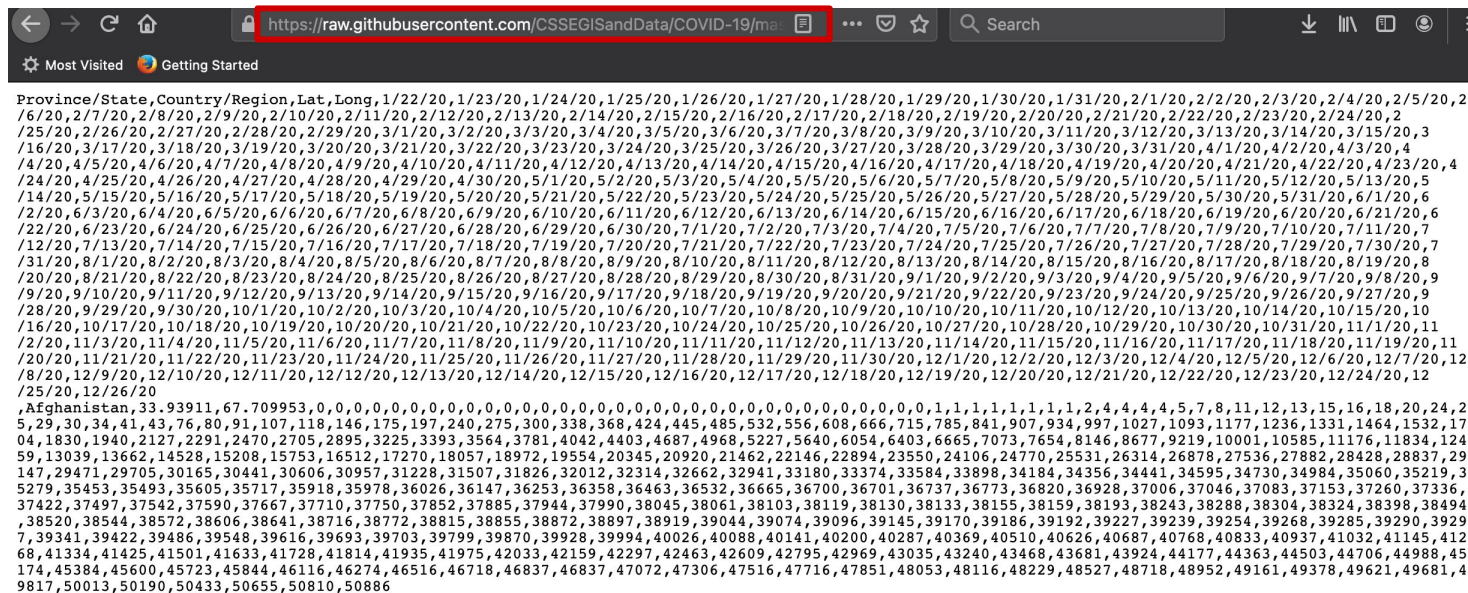
File Name	Description	Last Update
..		
.gitignore	update	11 months ago
Errata.csv	Update Errata.csv	2 days ago
README.md	Update README	8 months ago
time_series_covid19_confirmed_US.csv	Automated update	6 hours ago
time_series_covid19_confirmed_global.csv	Automated update	6 hours ago
time_series_covid19_deaths_US.csv	Automated update	6 hours ago
time_series_covid19_deaths_global.csv	Automated update	6 hours ago
time_series_covid19_recovered_global.csv	Automated update	6 hours ago

# Importing files from the internet

[illegible]

# Importing files from the internet

- Save the raw file as plain text file and download to your local machine
- Copy the link the of the raw file and download in R



# Importing files from the internet

```
# Require curl package
> install.packages("curl")
> library(curl)

# Create an object for the link
> url <-
"https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/cs
se_covid_19_data/csse_covid_19_time_series/time_series_covid19_confir
med_global.csv"

> covid_confirmed <- read_csv(url)
> dim(covid_confirmed)
```

[1] 271 343



# Importing files from the internet

```
> head(covid_confirmed)
```

```
# A tibble: 6 x 343
```

	`Province/State`	`Country/Region`	Lat	Long	`1/22/20`	`1/23/20`
	<chr>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
1	NA	Afghanistan	33.9	67.7	0	0
2	NA	Albania	41.2	20.2	0	0
3	NA	Algeria	28.0	1.66	0	0
4	NA	Andorra	42.5	1.52	0	0
5	NA	Angola	-11.2	17.9	0	0
6	NA	Antigua and Bar...	17.1	-61.8	0	0 ...

```
> tail(covid_confirmed)
```

```
# A tibble: 6 x 343
```

	`Province/State`	`Country/Region`	Lat	Long	`1/22/20`	`1/23/20`
	<chr>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
1	NA	Venezuela	6.42	-66.6	0	0
2	NA	Vietnam	14.1	108.	0	2
3	NA	West Bank and G...	32.0	35.2	0	0
4	NA	Yemen	15.6	48.5	0	0
5	NA	Zambia	-13.1	27.8	0	0
6	NA	Zimbabwe	-19.0	29.2	0	0 ...

# Writing to a file

**write\_csv**(x, file)

x = data frame or tibble

file = file to write into

```
> write_csv(covid_confirmed, file = "covid_confirmed.csv")
```

# Writing to a file

- **RDS** is R's custom binary format
- save column type specification

**write\_rds**(x, file)

x = data frame or tibble

file = file to write into

```
> write_rds(covid_confirmed, file = "covid_confirmed.rds")
```

```
# open RDS file
```

```
> read_rds("covid_confirmed.rds")
```

# Writing to a file

- **RDS** is R's custom binary format
- save column type specification

```
# open RDS file  
> read_rds("covid_confirmed.rds")
```

```
# A tibble: 271 x 343
```

	`Province/State` <chr>	`Country/Region` <chr>	Lat <dbl>	Long <dbl>	`1/22/20` <dbl>	`1/23/20` <dbl>
1	NA	Afghanistan	33.9	67.7	0	0
2	NA	Albania	41.2	20.2	0	0
3	NA	Algeria	28.0	1.66	0	0
4	NA	Andorra	42.5	1.52	0	0
5	NA	Angola	-11.2	17.9	0	0

```
...
```

# Take-away message

- readr package of tidyverse is useful for reading and writing textfiles
- readxl for importing Excel files in R
- use R to fetch files from the internet