Getting Data

Considerations for deciding how to get data

- reproducibility of workflow
- frequency with which data is updated
- available formats (may not be identical)
- time to process data

Web scraping

- Web scraping is a last resort, other methods are generally preferable if available
- Better to find an API, use httr package
- Even better, find an R package
 ex. https://cran.r-project.org/web/packages/atus/index.html

Case study: CDC birth data

Options:

- 1. .txt file from CDC https://www.cdc.gov/nchs/data access/vitalstatsonline.htm
- 2. .csv file from NBER https://www.nber.org/research/data/vital-statistics-natality-birth-data (2.46GB unzipped, 200MB zipped)
- 3. CDC Wonder API web interface https://wonder.cdc.gov/
- 4. CDC Wonder API https://github.com/socdataR/wonderapi

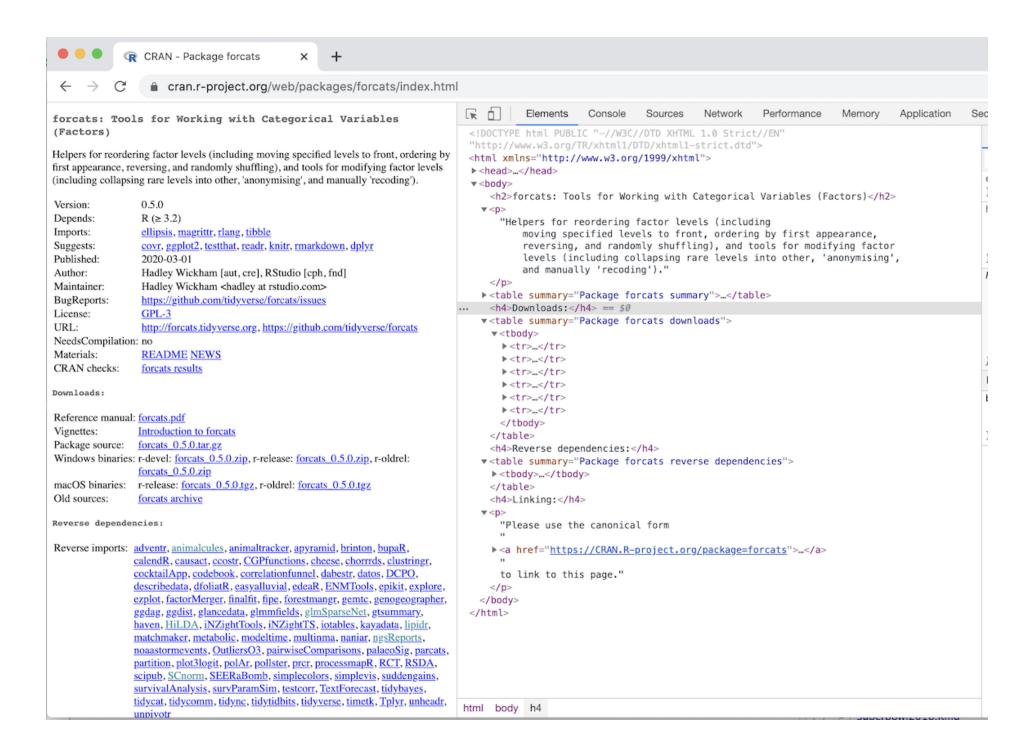
Web scraping, what not to do

- Scrape all Southwest Airlines data and send consumers notifications if their ticket prices decreased after purchase
- Buy a International Council of Shopping Centers membership, agree to terms of membership, then scrape the entire proprietary membership directory and contact members
- Scrape data that is for sale

Web scraping, what you should do

- think and investigate legal issues
- think about ethical questions
- limit bandwidth use
- scrape only what you need

Structure of an HTML page



rvest package

```
library(tidyverse)
library(rvest)
library(robotstxt)

paths_allowed("https://cran.r-project.org/web/packages/forcats/index.html")
```

[1] TRUE

paths_allowed("https://cran.r-project.org/web/packages/forcats/DESCRIPTION")

[1] FALSE

Tables

```
forcats_data <- read_html("https://cran.r-project.org/web/packages/forcats/index.html") %>%
  html_table()
length(forcats_data)
```

[1] 4

forcats_data[[1]]

X1 X2

Version: 0.5.2

Depends: $R (\geq 3.4)$

Imports: cli, ellipsis, glue, lifecycle, magrittr, rlang (≥

1.0.0), tibble, withr

Suggests: covr, dplyr, ggplot2, knitr, readr, rmarkdown,

testthat (≥

3.0.0)

Published: 2022-08-19

Author: Hadley Wickham [aut, cre],

RStudio [cph, fnd]

Maintainer: Hadley Wickham <hadley at rstudio.com>

X1 X2

BugReports: https://github.com/tidyverse/forcats/issues

License: MIT + file LICENSE

URL: https://forcats.tidyverse.org/,

https://github.com/tidyverse/forcats

NeedsCompilation: no

Materials: README NEWS

CRAN checks: forcats results

```
mytable <- forcats_data[[1]]
str(mytable)</pre>
```

```
## tibble [13 × 2] (S3: tbl_df/tbl/data.frame)
## $ X1: chr [1:13] "Version:" "Depends:" "Imports:" "Suggests:" ...
## $ X2: chr [1:13] "0.5.2" "R (\geq 3.4)" "cli, ellipsis, glue, lifecycle, magrittr, rlang (\geq 1.0.0),
tibble, withr" "covr, dplyr, ggplot2, knitr, readr, rmarkdown, testthat (\geq \n3.0.0)" ...
```

```
version <- mytable %>% filter(X1 == "Version:") %>% pull(X2)
date <- mytable %>% filter(X1 == "Published:") %>% pull(X2)
```

The most recent version of **forcats** on CRAN is 0.5.2, published on 2022-08-19.

(Use <u>inline rmarkdown syntax</u> to include values of variables within text sections.)

Data not in table form

https://www.beckershospitalreview.com/public-health/states-ranked-by-percentage-of-covid-19-vaccines-administered.html

```
vaccine <- read_html("https://www.beckershospitalreview.com/public-health/states-ranked-by-percentage-of-covid-19-vaccines-
administered.html")</pre>
```

```
vaccine |> html_element("#inner-article-content")
```

```
## {html node}
## <div id="inner-article-content">
         [1] Wisconsin has administered the highest percentage of COVID-19 vaccine ...
           [2] <script type="text/javascript">doNotShowRelatedArticles = 1;</script>
          [3] The <a href="https://covid.cdc.gov/covid-data-tracker/#vaccinations" ...
            [4] As of 6 a.m. ET Nov. 29, a total of 570,662,725 vaccine doses had bee ...
            [5] Below are the states and Washington, D.C., ranked by the percentage o ...
            [6] 1. <strong>Wisconsin</strong>Construction of the state of the s
           [7] 2. <strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecticut</strong>Connecti
          [8] 3. <strong>Massachusetts</strong><br>Doses distributed to state: 13,3 ...
         [9] 4. <strong>New Mexico</strong><br>Doses distributed to state: 3,599,3 ...
## [10] 5. <strong>Vermont</strong><br>Doses distributed to state: 1,295,970< ...
## [11] 6. <strong>Rhode Island</strong><br>Doses distributed to state: 2,020 ...
## [12] 7. <strong>Colorado</strong>Cobr>Doses distributed to state: 10,087,26 ...
## [13] 8. <strong>California</strong><br>Doses distributed to state: 70,222, ...
## [14] 9. <strong>New York State</strong><br>Doses distributed to state: 35, ...
## [15] 10. <strong>Virginia</strong><br>Doses distributed to state: 15,561,3 ...
## [16] 11. <strong>Maine</strong><br>Doses distributed to state: 2,642,860<br/>...
## [17] 12. <strong>Illinois</strong><br>Doses distributed to state: 21,451,2 ...
## [18] 13. <strong>Minnesota</strong><br>Doses distributed to state: 9,786,0 ...
## [19] 14. <strong>Nevada</strong><br>Doses distributed to state: 4,757,360< ...
## [20] 15. <strong>Arizona</strong><br>Doses distributed to state: 11,680,64 ...
## ...
```

Troubleshooting

- rvest makes it easy to identify nodes and parse text
- but... it doesn't work with all dynamically created content
- workaround: download page as "Webpage, complete" manually
- Or: use RSelenium

Example

https://analytics.usa.gov/

```
h2 tag

html_elements("h2")

id attribute

html_elements("#current_visitors")

class attribute

html_elements(".data")
```

Examples

```
library(robotstxt)
paths_allowed("https://analytics.usa.gov/")

## [1] TRUE

webdata <- read_html("https://analytics.usa.gov/")
webdata %>% html_elements("h2")

## {xml_nodeset (1)}
## [1] <h2 id="current_visitors" class="data">...</h2>

webdata %>% html_element("#current_visitors")
```

```
## {xml nodeset (16)}
  [1] <h2 id="current visitors" class="data">...</h2>
  [2] <svg class="data time-series"></svg>
  [3] <span id="total visitors" class="data">...</span>
   [4] <div class="data bar-chart">\n
                                                  </div>
   [5] <div class="data bar-chart">\n
                                                  </div>
   [6] <div class="data bar-chart">\n
                                                  </div>
   [7] <div class="data bar-chart">\n
                                                  </div>
  [8] <div class="data bar-chart">\n
                                                  </div>
  [9] <div class="data bar-chart">\n
                                                  </div>
## [10] <div class="data bar-chart">\n
                                                  </div>
## [11] <div class="data bar-chart">\n
                                                  </div>
## [12] <div class="data bar-chart">\n
                                                  </div>
## [13] <div class="data bar-chart">\n
                                                  </div>
## [14] <div class="data bar-chart">\n
                                                  </div>
## [15] <div class="data bar-chart">\n
                                                  </div>
## [16] <div class="data bar-chart">\n
                                                  </div>
```

```
webdata %>% html_elements("h2") %>% html_text()
```

```
## [1] "..."
```

Where's the number?

```
webdata_dl <- read_html("analytics.html")
webdata_dl %>% html_elements("h2") %>% html_text()
```

```
## [1] "414,029"
```

```
webdata dl %>% html elements(".data")
```

```
## {xml nodeset (16)}
  [1] <h2 id="current visitors" class="data">414,029</h2>
  [2] <svq class="data time-series" viewbox="0 0 700 150"><q class="axis y0" t ...
  [3] span id="total visitors" class="data">5.60 billion
   [4] <div class="data bar-chart">\n
                                                  <div class="bin">\n<div class= ...
  [5] <div class="data bar-chart">\n
                                                 <div class="bin">\n<div class= ...
  [6] <div class="data bar-chart">\n
                                                 <div class="bin">\n<div class= ...
  [7] <div class="data bar-chart">\n
                                                  <div class="bin">\n<div class= ...
   [8] <div class="data bar-chart">\n
                                                 <div class="bin" data-share="2 ...
  [9] <div class="data bar-chart">\n
                                                 <div class="bin">\n<div class= ...
## [10] <div class="data bar-chart">\n
                                                 <div class="bin">\n<div class= ...
## [11] <div class="data bar-chart">\n
                                                 <div class="bin" data-share="8 ...
                                                 <div class="bin" data-share="1 ...
## [12] <div class="data bar-chart">\n
## [13] <div class="data bar-chart">\n
                                                  <div class="bin">\n<div class= ...
## [14] <div class="data bar-chart">\n
                                                 <div class="bin">\n<div class= ...
                                                 <div class="bin">\n<div class= ...
## [15] <div class="data bar-chart">\n
## [16] <div class="data bar-chart">\n
                                                 <div class="bin">\n<div class= ...
```