Lecture 19: Web scraping and data wrangling

Motivation: Taskmaster



Last time:

- read_html to import web page
- html_elements to extract elements of interest
- html_table to extract tabular data
- html_text2 to extract text data

Series 11







```
Prev 1 • 2 • 3 • 4 • 5 • CoC • 6 • 7 • 8 • 9 • 10 • NYT • 11 • 12 • NYT || • 13 • CoC || • 14 • NYT ||| • 15
                                              • 16 Next
```

```
It's not your fault. • The Lure of the Treacle Puppies. • Run up a tree to the moon. •
Premature conker. • Slap and tong. • Absolute casserole. • You've got no chutzpah. • An
           orderly species. • Mr Octopus and Pottyhands. • Activate Jamali.
```

Series 11



```
Prev 1 • 2 • 3 • 4 • 5 • CoC • 6 • 7 • 8 • 9 • 10 • NYT • 11 • 12 • NYT III • 13 • CoC II • 14 • NYT III • 15 • 16 Next
```

It's not your fault. • The Lure of the Treacle Puppies. • Run up a tree to the moon. • Premature conker. • Slap and tong. • Absolute casserole. • You've got no chutzpah. • An orderly species. • Mr Octopus and Pottyhands. • Activate Jamali.

```
1 read html("https://taskmaster.fandom.com/wiki/Series 11") |>
     html elements("td")
{xml nodeset (436)}
 [1] \n<table class="toccolours" align="center" style="background:
#89110 ...
 [2] \n<a href="/wiki/Series 10" title="Series 10"><span
style="color: #F ...
 [3] \n<span style="font-family: Veteran
Typewriter; "><a h ...
 [4] <td class="pi-horizontal-group-item pi-data-value pi-font pi-
border-colo ...
 [5] <td class="pi-horizontal-group-item pi-data-value pi-font pi-
border-colo ...
 [6] Episode 1: <span style="font-family: Veteran"
Typewriter; ...
 [7] \n<a href="/wiki/Best thing you can carry, but only just"
```

```
1 read_html("https://taskmaster.fandom.com/wiki/Series_11") |>
2 html_elements("td[align='center']")

{xml_nodeset (1)}
[1] \n<span style="font-family: Veteran Typewriter;">
<a hr ...</pre>
```

```
1 read html("https://taskmaster.fandom.com/wiki/Series 11") |>
      html elements("td[align='center'] > span")
{xml nodeset (10)}
 [1] <span style="font-family: Veteran Typewriter;"><a
href="/wiki/It%27s not ...
 [2] span style="font-family: Veteran Typewriter;"><a</pre>
href="/wiki/The Lure o ...
 [3] <span style="font-family: Veteran Typewriter;"><a
href="/wiki/Run up a t ...
 [4] span style="font-family: Veteran Typewriter;"><a</pre>
href="/wiki/Premature ...
 [5] <span style="font-family: Veteran Typewriter;"><a
href="/wiki/Slap and t ...
 [6] <span style="font-family: Veteran Typewriter;"><a
href="/wiki/Absolute c ...
 [7] span style="font-family: Veteran Typewriter;"><a</pre>
```

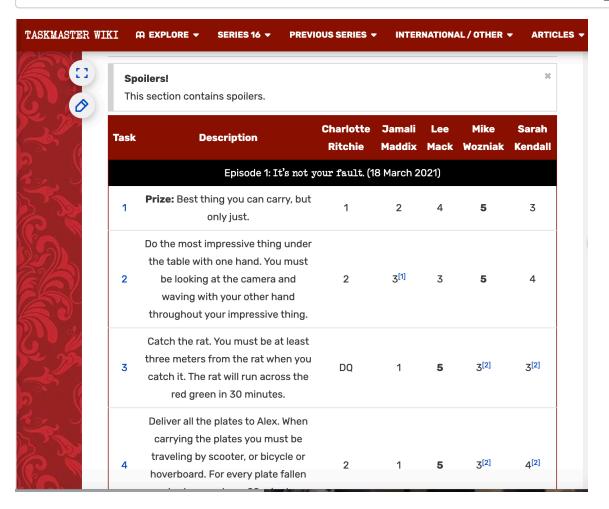
```
1 read html("https://taskmaster.fandom.com/wiki/Series 11") |>
      html elements("td[align='center'] > span > a")
{xml nodeset (10)}
 [1] <a href="/wiki/It%27s not your fault." title="It's not your
fault.">It's ...
 [2] <a href="/wiki/The Lure of the Treacle Puppies." title="The Lure of
the ...
 [3] <a href="/wiki/Run_up_a_tree_to_the_moon." title="Run up a tree to
the m ...
 [4] <a href="/wiki/Premature conker." title="Premature
conker.">Premature co ...
 [5] <a href="/wiki/Slap and tong." title="Slap and tong.">Slap and
tong.</a>
 [6] <a href="/wiki/Absolute casserole." title="Absolute
casserole.">Absolute ...
 [7] <a href="/wiki/You%27ve got no chutzpah." title="You've got no
```

```
1 read_html("https://taskmaster.fandom.com/wiki/Series_11") |>
2 html_elements("td[align='center'] > span > a") |>
3 html_attr("href")

[1] "/wiki/It%27s_not_your_fault."
[2] "/wiki/The_Lure_of_the_Treacle_Puppies."
[3] "/wiki/Run_up_a_tree_to_the_moon."
[4] "/wiki/Premature_conker."
[5] "/wiki/Slap_and_tong."
[6] "/wiki/Absolute_casserole."
[7] "/wiki/You%27ve_got_no_chutzpah."
[8] "/wiki/An_orderly_species."
[9] "/wiki/Mr_Octopus_and_Pottyhands."
[10] "/wiki/Activate_Jamali."
```

Last time: Taskmaster data

https://taskmaster.fandom.com/wiki/Series_11



Scraping the tabular data

```
read html("https://taskmaster.fandom.com/wiki/Series 11") |>
      html element(".tmtable") |>
      html table()
# A tibble: 75 \times 7
                        Description `Charlotte Ritchie` `Jamali Maddix`
   Task
`Lee Mack`
   <chr>
                        <chr>
                                     <chr>
                                                           <chr>
<chr>
 1 Episode 1: It's n... Episode 1:... Episode 1: It's no... Episode 1: It'...
Episode 1...
 2 1
                        Prize: Bes... 1
                                                            2
 3 2
                        Do the mos... 2
                                                            3[1]
 4 3
                        Catch the ... DQ
 5 4
                        Deliver al... 2
 6 5
                        Live: Stac... 0
                                                                             0
 7 Total
                        Total
                                                                             17
 8 Episode 2: The Lu... Episode 2:... Episode 2: The Lur... Episode 2: The...
```

Here's what we have so far:

```
Description `Charlotte Ritchie` `Jamali Maddix`
      Task
    Episode 1: It's... Episode 1:... Episode 1: It's no... Episode 1: It'... Epi
                        Prize: Bes... 1
                        Do the mos... 2
                                                             3[1]
                        Catch the ... DO
                        Deliver al... 2
                        Live: Stac... 0
    Total
                        Total
                                                                               17
    Episode 2: The ... Episode 2: ... Episode 2: The Lur ... Episode 2: The ... Epi
                        Prize: Bes... 5
10
11
                        Make the b... 0
                                                             5
                                                                               0
```

What changes do you think we should make to this format?

What we ultimately want:

```
Task
            Description
                               episode episode name air date contestant sco
                                        "It's not y... 18 Marc... Charlotte... 1
             Prize: Best th... 1
             Prize: Best th... 1
                                        "It's not y... 18 Marc... Jamali Ma... 2
             Prize: Best th... 1
                                        "It's not y... 18 Marc... Lee Mack
         Prize: Best th... 1
                                        "It's not y... 18 Marc... Mike Wozn... 5
                                        "It's not y... 18 Marc... Sarah Ken... 3
         Prize: Best th... 1
             Do the most im... 1
                                        "It's not y... 18 Marc... Charlotte... 2
         Do the most im... 1
                                        "It's not y... 18 Marc... Jamali Ma... 3
 9
             Do the most im... 1
                                        "It's not y... 18 Marc... Lee Mack
10
             Do the most im... 1
                                        "It's not y... 18 Marc... Mike Wozn... 5
   10 2
             Do the most im... 1
                                        "It's not y... 18 Marc... Sarah Ken... 4
```

colnames: Task, Description, episode, episode_name, air_date, contestant, score, series

Intermediate step:

```
Task Description
                                     episode
                                                contestant score series
             Prize: Best thing...
                                   Episode 1... Charlotte... 1
      1
                                                                         11
            Prize: Best thing...
                                   Episode 1... Jamali Ma... 2
                                                                         11
            Prize: Best thing...
                                   Episode 1... Lee Mack
                                                                         11
 4
 5
            Prize: Best thing...
                                   Episode 1... Mike Wozn... 5
                                                                         11
                                   Episode 1... Sarah Ken... 3
 6
            Prize: Best thing...
                                                                        11
            Do the most...
                                   Episode 1... Charlotte... 2
                                                                        11
                                   Episode 1... Jamali Ma... 3
 8
            Do the most...
                                                                        11
 9
      2
            Do the most...
                                   Episode 1... Lee Mack
                                                                        11
      2
                                   Episode 1... Mike Wozn... 5
10
             Do the most...
                                                                         11
11
                                   Episode 1... Sarah Ken... 4
      2
             Do the most...
                                                                        11
```

Here's what we have so far:

```
Description `Charlotte Ritchie` `Jamali Maddix`
      Task
    Episode 1: It's... Episode 1:... Episode 1: It's no... Episode 1: It'... Epi
                        Prize: Bes... 1
                        Do the mos... 2
                                                             3[1]
                        Catch the ... DO
                        Deliver al... 2
                        Live: Stac... 0
    Total
                        Total
                                                                               17
    Episode 2: The ... Episode 2: ... Episode 2: The Lur ... Episode 2: The ... Epi
                        Prize: Bes... 5
10
                        Make the b... 0
11
                                                                               0
```

What wrangling steps do we need to take?

Step 1: create a separate column for episode:

```
1 read html("https://taskmaster.fandom.com/wiki/Series 11") |>
      html element(".tmtable") |>
      html table() |>
      mutate(episode = ifelse(startsWith(Task, "Episode"), Task, NA))
# A tibble: 75 \times 2
                                                                   episode
   Task
   <chr>
                                                                   <chr>
 1 Episode 1: It's not your fault. (18 March 2021)
                                                                   Episode
1: It's ...
 2 1
                                                                   <NA>
 3 2
                                                                   < NA >
 4 3
                                                                   < NA >
 5 4
                                                                   <NA>
 6 5
                                                                   <NA>
 7 Total
                                                                   <NA>
 8 Episode 2: The Lure of the Treacle Puppies. (25 March 2021) Episode
2: The L...
 9 1
                                                                   < NA >
```

Step 2: fill in the episodes

```
read html("https://taskmaster.fandom.com/wiki/Series 11") |>
     html element(".tmtable") |>
     html table() |>
     mutate(episode = ifelse(startsWith(Task, "Episode"), Task, NA)) |>
 4
     fill(episode, .direction = "down")
      Task
                                         episode
      Episode 1: It's...
                                       Episode 1: It'...
 3
                                        Episode 1: It'...
                                        Episode 1: It'...
      2
 4
                                        Episode 1: It'...
 6
      4
                                        Episode 1: It'...
      5
                                        Episode 1: It'...
 8
      Total
                                        Episode 1: It'...
      Episode 2: The Lure of...
                                       Episode 2: The...
                                        Episode 2: The...
10
      1
                                        Episode 2: The...
11
      2
```

Step 3: remove the "Total" and "Episode" rows in the Task column

```
Task
                                           episode
      Episode 1: It's...
                                         Episode 1: It'...
                                         Episode 1: It'...
                                         Episode 1: It'...
 4
                                         Episode 1: It'...
 6
                                         Episode 1: It'...
       5
                                         Episode 1: It'...
       Total
                                         Episode 1: It'...
                                         Episode 2: The...
       Episode 2: The Lure of...
                                         Episode 2: The...
10
       1
                                         Episode 2: The...
11
```

What R code would we use to remove these rows?

Step 3: remove the "Total" and "Episode" rows in the Task column

```
read html("https://taskmaster.fandom.com/wiki/Series 11") |>
      html element(".tmtable") |>
      html table() |>
      mutate(episode = ifelse(startsWith(Task, "Episode"), Task, NA)) |>
      fill(episode, .direction = "down") |>
      filter(!startsWith(Task, "Episode"),
             !(Task %in% c("Total", "Grand Total")))
# A tibble: 54 \times 8
                                   `Charlotte Ritchie` `Jamali Maddix`
   Task Description
`Lee Mack`
                                                       <chr>
   <chr> <chr>
                                   <chr>
<chr>
 1 1
         Prize: Best thing you c... 1
 2 2
         Do the most impressive ... 2
                                                       3[1]
 3 3
         Catch the rat. You must... DQ
 4 4
        Deliver all the plates ... 2
 5 5
         Live: Stack your bucket... 0
                                                                        0
         Prize: Best drinking ve... 5
 6 1
 7 2
         Make the balloon hover ... 0
                                                                        0
 8 3
         Team: Have an argument... 2
                                                                        5
```

	 	-	-	-

Step 4: Pivot

```
# A tibble: 54 \times 8
                                   `Charlotte Ritchie` `Jamali Maddix`
   Task Description
`Lee Mack`
   <chr> <chr>
                                   <chr>
                                                        <chr>
<chr>
 1 1
         Prize: Best thing you c... 1
 2 2
         Do the most impressive ... 2
                                                        3[1]
 3 3
         Catch the rat. You must... DO
 4 4 Deliver all the plates ... 2
         Live: Stack your bucket... 0
 6 1
         Prize: Best drinking ve... 5
         Make the balloon hover ... 0
                                                                         0
 8 3
         Team: Have an argument... 2
 9 4
         Make the house haunted. 3
```

How should we pivot this data?

Step 4: Pivot

```
read html("https://taskmaster.fandom.com/wiki/Series 11") |>
      html element(".tmtable") |>
      html table() |>
      mutate(episode = ifelse(startsWith(Task, "Episode"), Task, NA)) |>
      fill(episode, .direction = "down") |>
      filter(!startsWith(Task, "Episode"),
 6
              !(Task %in% c("Total", "Grand Total"))) |>
      pivot longer(cols = -c(Task, Description, episode),
 9
                    names to = "contestant",
                    values to = "score") |>
10
11
      mutate(series = 11)
# A tibble: 270 \times 6
   Task Description
                                                   episode contestant
score series
   <chr> <chr>
                                                   <chr>
                                                           <chr>
<chr> <dbl>
         Prize: Best thing you can carry, but o... Episod... Charlotte... 1
 1 1
11
 2 1
         Prize: Best thing you can carry, but o... Episod... Jamali Ma... 2
11
3 1
         Prize: Best thing you can carry, but o... Episod... Lee Mack
11
```

Prize: Best thing you can carry, but o... Episod... Mike Wozn... 5

Prize: Best thing you can carry but o Episod Sarah Ken 3

Next steps

- Separate episode info into episode number, episode name, and air date columns
- Combine data from multiple series