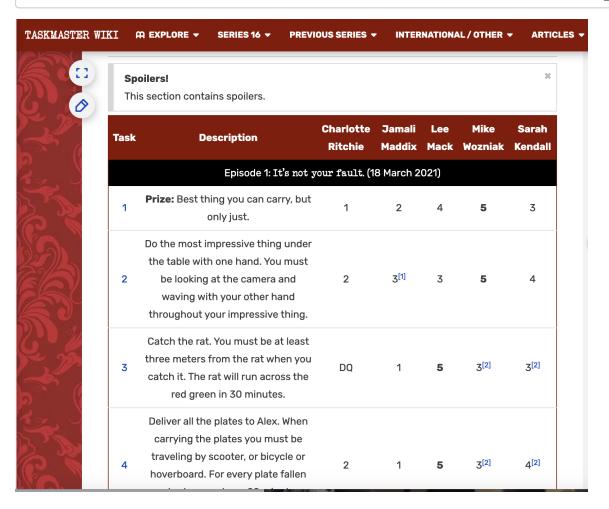
# Lecture 20: Web scraping and data wrangling

## Last time: Taskmaster data

1 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
             Prize: Best th... 1
                                       "It's not y... 18 Marc... Charlotte... 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[1
 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
                                        "It's not y... 18 Marc... Sarah Ken... 4
   10 2
             Do the most im... 1
```

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