

# RWorksheet#5

Barrientos, Delfin, Infiesto

2024-11-06

#Extracting TV Shows Reviews

*#1. Each group needs to extract the top 50 tv shows in Imdb.com. It will include the rank, the title of  
#tv show, tv rating, the number of people who voted, the number of episodes, the year it was released.  
#It will also include the number of user reviews and the number of critic reviews, as well as the popul  
#rating for each tv shows.*

```
library(dplyr)
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
## filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
## intersect, setdiff, setequal, union
```

```
library(rvest)
```

```
library(polite)
```

```
library(tidyr)
```

```
library(httr)
```

```
url <- "https://www.imdb.com/chart/toptv/?ref_=nv_tvv_250"
```

```
session <- bow(url, user_agent = "Educational Purposes Only")
```

```
page <- scrape(session)
```

```
titles <- page %>%
```

```
  html_nodes('.ipc-title__text') %>%
```

```
  html_text()
```

```
ranks <- page %>%
```

```
  html_nodes('a.ipc-title-link-wrapper') %>%
```

```

html_text() %>%
gsub("[^0-9]", "", .)

links <- page %>%
  html_nodes('a.ipc-title-link-wrapper') %>%
  html_attr('href') %>%
  paste0("https://www.imdb.com", .)

titles <- titles[!titles %in% c("IMDb Charts", "Top 250 TV Shows")]
titles <- titles[1:50]
ranks <- ranks[1:50]
links <- links[1:50]

rank_title <- data.frame(
  rank = ranks,
  title = titles,
  link = links,
  stringsAsFactors = FALSE
)

scrape_show_details <- function(show_url) {
  Sys.sleep(1) # Be polite, avoid overwhelming IMDb's servers
  page <- tryCatch(read_html(show_url), error = function(e) return(NULL))

  if (is.null(page)) return(NULL)

  rating <- page %>%
    html_node('.sc-d541859f-1 imUuxf') %>%
    html_text(trim = TRUE)

  votes <- page %>%
    html_node('div[class*="AggregateRatingButton__TotalRatingAmount"]') %>%
    html_text(trim = TRUE)

  episodes <- page %>%
    html_node('a[href*="episodes?"]') %>%
    html_text(trim = TRUE)

  year <- page %>%
    html_node('span[class*="TitleBlockMetaData__ListItemText"]') %>%
    html_text(trim = TRUE)

  user_reviews <- page %>%
    html_node('span[data-testid="reviews-header"]') %>%
    html_text(trim = TRUE)

  critic_reviews <- page %>%
    html_node('span[class*="score"]') %>%
    html_text(trim = TRUE)

```

```

popularity <- page %>%
  html_node('div[class*="trending-list-rank-item__position"]') %>%
  html_text(trim = TRUE)

return(data.frame(
  rating = ifelse(is.na(rating), "N/A", rating),
  votes = ifelse(is.na(votes), "N/A", votes),
  episodes = ifelse(is.na(episodes), "N/A", episodes),
  year = ifelse(is.na(year), "N/A", year),
  user_reviews = ifelse(is.na(user_reviews), "N/A", user_reviews),
  critic_reviews = ifelse(is.na(critic_reviews), "N/A", critic_reviews),
  popularity = ifelse(is.na(popularity), "N/A", popularity),
  stringsAsFactors = FALSE
))
}

show_details <- lapply(rank_title$link, function(link_url) {
  scrape_show_details(link_url)
})

show_details <- Filter(Negate(is.null), show_details)
final_data <- do.call(rbind, show_details)

final_result <- cbind(rank_title, final_data)

write.csv(final_result, file = "top_50_tv_shows_imdb.csv", row.names = FALSE)

final_result

```

```

##      rank      title
## 1      1  1. Breaking Bad
## 2      2  2. Planet Earth II
## 3      3  3. Planet Earth
## 4      4  4. Band of Brothers
## 5      5  5. Chernobyl
## 6      6  6. The Wire
## 7      7  7. Avatar: The Last Airbender
## 8      8  8. Blue Planet II
## 9      9  9. The Sopranos
## 10     10 10. Cosmos: A Spacetime Odyssey
## 11     11      11. Cosmos
## 12     12      12. Our Planet
## 13     13      13. Game of Thrones
## 14     14      14. Bluey
## 15     15      15. The World at War
## 16     16 16. Fullmetal Alchemist Brotherhood
## 17     17      17. Rick and Morty
## 18     18      18. Life

```

## 19	19	19. The Last Dance
## 20	20	20. The Twilight Zone
## 21	21	21. The Vietnam War
## 22	22	22. Sherlock
## 23	23	23. Attack on Titan
## 24	24	24. Batman: The Animated Series
## 25	25	25. The Office
## 26	<NA>	Recently viewed
## 27	<NA>	
## 28	<NA>	
## 29	<NA>	
## 30	<NA>	
## 31	<NA>	
## 32	<NA>	
## 33	<NA>	
## 34	<NA>	
## 35	<NA>	
## 36	<NA>	
## 37	<NA>	
## 38	<NA>	
## 39	<NA>	
## 40	<NA>	
## 41	<NA>	
## 42	<NA>	
## 43	<NA>	
## 44	<NA>	
## 45	<NA>	
## 46	<NA>	
## 47	<NA>	
## 48	<NA>	
## 49	<NA>	
## 50	<NA>	

##	link	rating	votes
## 1	<a href="https://www.imdb.com/title/tt0903747/?ref_=chttvtp_t_1">https://www.imdb.com/title/tt0903747/?ref_=chttvtp_t_1</a>	N/A	N/A
## 2	<a href="https://www.imdb.com/title/tt5491994/?ref_=chttvtp_t_2">https://www.imdb.com/title/tt5491994/?ref_=chttvtp_t_2</a>	N/A	N/A
## 3	<a href="https://www.imdb.com/title/tt0795176/?ref_=chttvtp_t_3">https://www.imdb.com/title/tt0795176/?ref_=chttvtp_t_3</a>	N/A	N/A
## 4	<a href="https://www.imdb.com/title/tt0185906/?ref_=chttvtp_t_4">https://www.imdb.com/title/tt0185906/?ref_=chttvtp_t_4</a>	N/A	N/A
## 5	<a href="https://www.imdb.com/title/tt7366338/?ref_=chttvtp_t_5">https://www.imdb.com/title/tt7366338/?ref_=chttvtp_t_5</a>	N/A	N/A
## 6	<a href="https://www.imdb.com/title/tt0306414/?ref_=chttvtp_t_6">https://www.imdb.com/title/tt0306414/?ref_=chttvtp_t_6</a>	N/A	N/A
## 7	<a href="https://www.imdb.com/title/tt0417299/?ref_=chttvtp_t_7">https://www.imdb.com/title/tt0417299/?ref_=chttvtp_t_7</a>	N/A	N/A
## 8	<a href="https://www.imdb.com/title/tt6769208/?ref_=chttvtp_t_8">https://www.imdb.com/title/tt6769208/?ref_=chttvtp_t_8</a>	N/A	N/A
## 9	<a href="https://www.imdb.com/title/tt0141842/?ref_=chttvtp_t_9">https://www.imdb.com/title/tt0141842/?ref_=chttvtp_t_9</a>	N/A	N/A
## 10	<a href="https://www.imdb.com/title/tt2395695/?ref_=chttvtp_t_10">https://www.imdb.com/title/tt2395695/?ref_=chttvtp_t_10</a>	N/A	N/A
## 11	<a href="https://www.imdb.com/title/tt0081846/?ref_=chttvtp_t_11">https://www.imdb.com/title/tt0081846/?ref_=chttvtp_t_11</a>	N/A	N/A
## 12	<a href="https://www.imdb.com/title/tt9253866/?ref_=chttvtp_t_12">https://www.imdb.com/title/tt9253866/?ref_=chttvtp_t_12</a>	N/A	N/A
## 13	<a href="https://www.imdb.com/title/tt0944947/?ref_=chttvtp_t_13">https://www.imdb.com/title/tt0944947/?ref_=chttvtp_t_13</a>	N/A	N/A
## 14	<a href="https://www.imdb.com/title/tt7678620/?ref_=chttvtp_t_14">https://www.imdb.com/title/tt7678620/?ref_=chttvtp_t_14</a>	N/A	N/A
## 15	<a href="https://www.imdb.com/title/tt0071075/?ref_=chttvtp_t_15">https://www.imdb.com/title/tt0071075/?ref_=chttvtp_t_15</a>	N/A	N/A
## 16	<a href="https://www.imdb.com/title/tt1355642/?ref_=chttvtp_t_16">https://www.imdb.com/title/tt1355642/?ref_=chttvtp_t_16</a>	N/A	N/A
## 17	<a href="https://www.imdb.com/title/tt2861424/?ref_=chttvtp_t_17">https://www.imdb.com/title/tt2861424/?ref_=chttvtp_t_17</a>	N/A	N/A
## 18	<a href="https://www.imdb.com/title/tt1533395/?ref_=chttvtp_t_18">https://www.imdb.com/title/tt1533395/?ref_=chttvtp_t_18</a>	N/A	N/A
## 19	<a href="https://www.imdb.com/title/tt8420184/?ref_=chttvtp_t_19">https://www.imdb.com/title/tt8420184/?ref_=chttvtp_t_19</a>	N/A	N/A
## 20	<a href="https://www.imdb.com/title/tt0052520/?ref_=chttvtp_t_20">https://www.imdb.com/title/tt0052520/?ref_=chttvtp_t_20</a>	N/A	N/A
## 21	<a href="https://www.imdb.com/title/tt1877514/?ref_=chttvtp_t_21">https://www.imdb.com/title/tt1877514/?ref_=chttvtp_t_21</a>	N/A	N/A

## 22	<a href="https://www.imdb.com/title/tt1475582/?ref_=chttvtp_t_22">https://www.imdb.com/title/tt1475582/?ref_=chttvtp_t_22</a>	N/A	N/A
## 23	<a href="https://www.imdb.com/title/tt2560140/?ref_=chttvtp_t_23">https://www.imdb.com/title/tt2560140/?ref_=chttvtp_t_23</a>	N/A	N/A
## 24	<a href="https://www.imdb.com/title/tt0103359/?ref_=chttvtp_t_24">https://www.imdb.com/title/tt0103359/?ref_=chttvtp_t_24</a>	N/A	N/A
## 25	<a href="https://www.imdb.com/title/tt0386676/?ref_=chttvtp_t_25">https://www.imdb.com/title/tt0386676/?ref_=chttvtp_t_25</a>	N/A	N/A
## 26	<NA>	N/A	N/A
## 27	<NA>	N/A	N/A
## 28	<NA>	N/A	N/A
## 29	<NA>	N/A	N/A
## 30	<NA>	N/A	N/A
## 31	<NA>	N/A	N/A
## 32	<NA>	N/A	N/A
## 33	<NA>	N/A	N/A
## 34	<NA>	N/A	N/A
## 35	<NA>	N/A	N/A
## 36	<NA>	N/A	N/A
## 37	<NA>	N/A	N/A
## 38	<NA>	N/A	N/A
## 39	<NA>	N/A	N/A
## 40	<NA>	N/A	N/A
## 41	<NA>	N/A	N/A
## 42	<NA>	N/A	N/A
## 43	<NA>	N/A	N/A
## 44	<NA>	N/A	N/A
## 45	<NA>	N/A	N/A
## 46	<NA>	N/A	N/A
## 47	<NA>	N/A	N/A
## 48	<NA>	N/A	N/A
## 49	<NA>	N/A	N/A
## 50	<NA>	N/A	N/A
##	episodes year user_reviews critic_reviews popularity		
## 1	Episodes62 N/A	N/A	5K N/A
## 2	Episodes6 N/A	N/A	158 N/A
## 3	Episodes11 N/A	N/A	111 N/A
## 4	Episodes10 N/A	N/A	1K N/A
## 5	Episodes5 N/A	N/A	3.5K N/A
## 6	Episodes60 N/A	N/A	785 N/A
## 7	Episodes62 N/A	N/A	997 N/A
## 8	Episodes7 N/A	N/A	53 N/A
## 9	Episodes86 N/A	N/A	959 N/A
## 10	Episodes13 N/A	N/A	205 N/A
## 11	Episodes13 N/A	N/A	80 N/A
## 12	Episodes12 N/A	N/A	245 N/A
## 13	Episodes74 N/A	N/A	5.8K N/A
## 14	Episodes194 N/A	N/A	366 N/A
## 15	Episodes26 N/A	N/A	126 N/A
## 16	Episodes68 N/A	N/A	463 N/A
## 17	Episodes78 N/A	N/A	908 N/A
## 18	Episodes11 N/A	N/A	12 N/A
## 19	Episodes10 N/A	N/A	541 N/A
## 20	Episodes156 N/A	N/A	213 N/A
## 21	Episodes10 N/A	N/A	176 N/A
## 22	Episodes15 N/A	N/A	1K N/A
## 23	Episodes98 N/A	N/A	2.3K N/A
## 24	Episodes85 N/A	N/A	218 N/A

## 25	Episodes188	N/A	N/A	1.7K	N/A
## 26	Episodes62	N/A	N/A	5K	N/A
## 27	Episodes6	N/A	N/A	158	N/A
## 28	Episodes11	N/A	N/A	111	N/A
## 29	Episodes10	N/A	N/A	1K	N/A
## 30	Episodes5	N/A	N/A	3.5K	N/A
## 31	Episodes60	N/A	N/A	785	N/A
## 32	Episodes62	N/A	N/A	997	N/A
## 33	Episodes7	N/A	N/A	53	N/A
## 34	Episodes86	N/A	N/A	959	N/A
## 35	Episodes13	N/A	N/A	205	N/A
## 36	Episodes13	N/A	N/A	80	N/A
## 37	Episodes12	N/A	N/A	245	N/A
## 38	Episodes74	N/A	N/A	5.8K	N/A
## 39	Episodes194	N/A	N/A	366	N/A
## 40	Episodes26	N/A	N/A	126	N/A
## 41	Episodes68	N/A	N/A	463	N/A
## 42	Episodes78	N/A	N/A	908	N/A
## 43	Episodes11	N/A	N/A	12	N/A
## 44	Episodes10	N/A	N/A	541	N/A
## 45	Episodes156	N/A	N/A	213	N/A
## 46	Episodes10	N/A	N/A	176	N/A
## 47	Episodes15	N/A	N/A	1K	N/A
## 48	Episodes98	N/A	N/A	2.3K	N/A
## 49	Episodes85	N/A	N/A	218	N/A
## 50	Episodes188	N/A	N/A	1.7K	N/A

[View\(final\\_result\)](#)