

# Mini Project - IMDB Web Scraping

```
library(tidyverse)
library(rvest) #scrape data from internet
```

```
url <- "https://www.imdb.com/search/title/?groups=top_100&sort=user_rating,desc"
```

```
print(url)
```

```
[1] "https://www.imdb.com/search/title/?groups=top_100&sort=user_rating,desc"
```

```
imdb <- read_html(url)
```

```
imdb
```

```
{html_document}
<html xmlns:og="http://ogp.me/ns#" xmlns:fb="http://www.facebook.com/2008/fbml"
[1] <head>\n<meta http-equiv="Content-Type" content="text/html; charset=UTF-8 .
[2] <body id="styleguide-v2" class="fixed">\n                <img height="1" widt .
```

```
titles <- imdb %>%
  html_nodes("h3.lister-item-header")%>%
  html_text2()
```

```
title[1:10]
```

```
'1. The Shawshank Redemption (1994)' · '2. The Godfather (1972)' · '3. The Dark Knight (2008)' ·  
'4. The Lord of the Rings: The Return of the King (2003)' · '5. Schindler\'s List (1993)' ·  
'6. The Godfather Part II (1974)' · '7. 12 Angry Men (1957)' · '8. Pulp Fiction (1994)' · '9. Inception (2010)' ·  
'10. The Lord of the Rings: The Two Towers (2002)'
```

```
ratings <- imdb %>%  
  html_nodes("div.inline-block.ratings-imdb-rating") %>%  
  html_text2() %>%  
  as.numeric()
```

```
num_votes <- imdb %>%  
  html_nodes("p.sort-num_votes-visible") %>%  
  html_text2()
```

```
#build a dataset  
df <- data.frame(  
  title = titles,  
  rating = ratings,  
  num_vote = num_votes  
)  
df
```

A data.frame: 50 × 3

title	rating	num_vote
<chr>	<dbl>	<chr>
1. The Shawshank Redemption (1994)	9.3	Votes: 2,659,313   Gross: \$28.34M   Top 250: #1
2. The Godfather (1972)	9.2	Votes: 1,842,945   Gross: \$134.97M   Top 250: #2
3. The Dark Knight (2008)	9.0	Votes: 2,632,089   Gross: \$534.86M   Top 250: #3
4. The Lord of the Rings: The Return of the King (2003)	9.0	Votes: 1,833,622   Gross: \$377.85M   Top 250: #7
5. Schindler's List (1993)	9.0	Votes: 1,346,943   Gross: \$96.90M   Top 250: #6
6. The Godfather Part II (1974)	9.0	Votes: 1,262,572   Gross: \$57.30M   Top 250: #4
7. 12 Angry Men (1957)	9.0	Votes: 785,202   Gross: \$4.36M   Top 250: #5
8. Pulp Fiction (1994)	8.9	Votes: 2,035,398   Gross: \$107.93M   Top 250: #8
9. Inception (2010)	8.8	Votes: 2,332,345   Gross: \$292.58M   Top 250: #14
10. The Lord of the Rings: The Two Towers (2002)	8.8	Votes: 1,655,684   Gross: \$342.55M   Top 250: #13
11. Fight Club (1999)	8.8	Votes: 2,104,435   Gross: \$37.03M   Top 250: #12
12. The Lord of the Rings: The Fellowship of the Ring (2001)	8.8	Votes: 1,862,422   Gross: \$315.54M   Top 250: #9
13. Forrest Gump (1994)	8.8	Votes: 2,060,490   Gross: \$330.25M   Top 250: #11
14. Il buono, il brutto, il cattivo (1966)	8.8	Votes: 758,155   Gross: \$6.10M   Top 250: #10
15. The Matrix (1999)	8.7	Votes: 1,900,520   Gross: \$171.48M   Top 250: #16
16. Goodfellas (1990)	8.7	Votes: 1,152,247   Gross: \$46.84M   Top 250: #17
17. The Empire Strikes Back (1980)	8.7	Votes: 1,284,476   Gross: \$290.48M   Top 250: #15
18. One Flew Over the Cuckoo's Nest (1975)	8.7	Votes: 1,003,271   Gross: \$112.00M   Top 250: #18
19. Interstellar (2014)	8.6	Votes: 1,804,741   Gross: \$188.02M   Top 250: #26
20. Cidade de Deus (2002)	8.6	Votes: 753,982   Gross: \$7.56M   Top 250: #23
21. Sen to Chihiro no kamikakushi (2001)	8.6	Votes: 757,260   Gross: \$10.06M   Top 250: #31
22. Saving Private Ryan (1998)	8.6	Votes: 1,382,090   Gross: \$216.54M   Top 250: #24
23. The Green Mile (1999)	8.6	Votes: 1,292,142   Gross: \$136.80M   Top 250: #27
24. La vita è bella (1997)	8.6	Votes: 691,486   Gross: \$57.60M   Top 250: #25
25. Se7en (1995)	8.6	Votes: 1,638,615   Gross: \$100.13M   Top 250: #19
26. Terminator 2: Judgment Day (1991)	8.6	Votes: 1,092,953   Gross: \$204.84M   Top 250: #29
27. The Silence of the Lambs (1991)	8.6	Votes: 1,422,638   Gross: \$130.74M   Top 250: #22
28. Star Wars (1977)	8.6	Votes: 1,357,110   Gross: \$322.74M   Top 250: #28

29. Seppuku (1962)	8.6	Votes: 57,087   Top 250: #44
30. Shichinin no samurai (1954)	8.6	Votes: 345,290   Gross: \$0.27M   Top 250: #20
31. It's a Wonderful Life (1946)	8.6	Votes: 453,771   Top 250: #21
32. Gisaengchung (2019)	8.5	Votes: 789,234   Gross: \$53.37M   Top 250: #34
33. Whiplash (2014)	8.5	Votes: 849,723   Gross: \$13.09M   Top 250: #42
34. The Intouchables (2011)	8.5	Votes: 852,739   Gross: \$13.18M   Top 250: #45
35. The Prestige (2006)	8.5	Votes: 1,324,835   Gross: \$53.09M   Top 250: #41
36. The Departed (2006)	8.5	Votes: 1,317,062   Gross: \$132.38M   Top 250: #39
37. The Pianist (2002)	8.5	Votes: 826,715   Gross: \$32.57M   Top 250: #33
38. Gladiator (2000)	8.5	Votes: 1,490,454   Gross: \$187.71M   Top 250: #37
39. American History X (1998)	8.5	Votes: 1,117,559   Gross: \$6.72M   Top 250: #38
40. The Usual Suspects (1995)	8.5	Votes: 1,080,443   Gross: \$23.34M   Top 250: #40
41. Léon (1994)	8.5	Votes: 1,154,058   Gross: \$19.50M   Top 250: #35
42. The Lion King (1994)	8.5	Votes: 1,051,565   Gross: \$422.78M   Top 250: #36
43. Nuovo Cinema Paradiso (1988)	8.5	Votes: 260,664   Gross: \$11.99M   Top 250: #52
44. Hotaru no haka (1988)	8.5	Votes: 276,128   Top 250: #46
45. Back to the Future (1985)	8.5	Votes: 1,195,749   Gross: \$210.61M   Top 250: #30
46. Apocalypse Now (1979)	8.5	Votes: 664,957   Gross: \$83.47M   Top 250: #53
47. Alien (1979)	8.5	Votes: 877,781   Gross: \$78.90M   Top 250: #50
48. Once Upon a Time in the West (1968)	8.5	Votes: 328,849   Gross: \$5.32M   Top 250: #48
49. Psycho (1960)	8.5	Votes: 669,353   Gross: \$32.00M   Top 250: #32
50. Rear Window (1954)	8.5	Votes: 490,491   Gross: \$36.76M   Top 250: #49

## Mini Project 02 - Specphone phone Database

```
library(tidyverse)
library(rvest) #scrape data from internet
```

```
url <- read_html("https://specphone.com/Samsung-Galaxy-A04.html")
```

```
att <- url %>%
  html_nodes("div.topic") %>%
  html_text2()

value <- url %>%
  html_nodes("div.detail") %>%
  html_text2()

df <- data.frame(
  atts = att,
  values = value
)
df
```

A data.frame: 31 × 2

atts	values
<chr>	<chr>
วันเปิดตัว	ตุลาคม 2565
วันวางจำหน่าย	ยังไม่วางจำหน่าย
ขนาด	164.40 x 76.30 x 9.10 มม.
น้ำหนัก	192 กรัม
วัสดุ	Glass front, plastic back, plastic frame
SIM	รองรับ 2 ซิมการ์ด (nano sim, nano sim)
Technology	HSPA 42.2/5.76 Mbps, LTE-A
2G	850/900/1800/1900
3G	850/900/1900/2100
4G	850/900/1900/2100/2600
5G	-
ความเร็ว	HSPA 42.2/5.76 Mbps, LTE-A
ประเภท	PLS LCD
ขนาดหน้าจอ	6.50 นิ้ว
ความละเอียด	720 x 1600 pixels
ระบบปฏิบัติการ	Android 12
ชิปประมวลผล	Spreadtrum Unisoc SC9863A 1.6 GHz
ชิปกราฟิก	PowerVR GE8322
หน่วยความจำ	3 GB
ความจุ	32 GB
Memory Card	microSD (1)
กล้องหลัก	ตัวที่ 1: 50 MP, f/1.8, (wide), AF ตัวที่ 2: 2 MP, f/2.4, (depth)
ความละเอียดวิดีโอ	1080p@30fps
กล้องหน้า	ตัวที่ 1: 5 MP, f/2.2
Bluetooth	5.0, A2DP, LE
Wi-Fi	802.11 a/b/g/n/ac, dual-b
USB	Type-C
GPS	GLONASS, GALILEO, BDS
NFC	ไม่รองรับ
ความจุ	5,000 mAh
ประเภท	Non-removable Li-Po Batt

```
#all samsung smartphone
samsung_url <- read_html("https://specphone.com/brand/Samsung")
```

```
#links to all samsung smartphone
links <- samsung_url %>%
  html_nodes("li.mobile-brand-item a") %>%
  html_attr("href")
links
```

```

'/Samsung-Galaxy-M13.html' · '/Samsung-Galaxy-A23.html' · '/Samsung-Galaxy-A13.html' ·
'/Samsung-Galaxy-M32-5G.html' · '/Samsung-Galaxy-A12-Nacho.html' · '/Samsung-Galaxy-Pocket-Neo.html' ·
'/Samsung-Galaxy-Young.html' · '/Samsung-Galaxy-J1-Mini.html' · '/Samsung-Galaxy-A01-Core-1-16GB.html' ·
'/Samsung-Galaxy-V-PLUS.html' · '/Samsung-Galaxy-Young-2.html' · '/Samsung-Galaxy-M02.html' ·
'/Samsung-Galaxy-A11.html' · '/Samsung-Galaxy-J2-Pro-2018.html' · '/Samsung-Galaxy-A12-2021.html' ·
'/Samsung-Galaxy-A21s-3-32GB.html' · '/Samsung-Galaxy-J5.html' · '/Samsung-Galaxy-J4.html' ·
'/Samsung-Galaxy-Core-2-Duos.html' · '/Samsung-Galaxy-Ace-Plus.html' · '/Samsung-Galaxy-A20.html' ·
'/Samsung-Galaxy-Chat.html' · '/Samsung-Galaxy-Gio.html' · '/Samsung-Galaxy-Tab-A7-Lite-LTE.html' ·
'/Samsung-Galaxy-Tab-A-10.5WIFI.html' · '/Samsung-Galaxy-Alpha.html' · '/Samsung-Galaxy-S3-Slim.html' ·
'/Samsung-Galaxy-S4-zoom.html' · '/Samsung-Galaxy-Xcover-2.html' · '/Samsung-Galaxy-Tab-8.9-3G-16GB.html' ·
'/Samsung-Galaxy-Tab-A8-LTE-2021.html' · '/Samsung-Galaxy-A8-2018.html' ·
'/Samsung-Galaxy-Tab4-8.0-wifi.html' · '/Samsung-Galaxy-M33-5G.html' · '/Samsung-Galaxy-A50.html' ·
'/Samsung-Galaxy-E7.html' · '/Samsung-Galaxy-S6.html' · '/Samsung-Galaxy-S20-FE.html' ·
'/Samsung-Galaxy-Tab-S4-WIFI.html' · '/Samsung-Galaxy-S7.html' · '/Samsung-Galaxy-Note-5-Exynos.html' ·
'/Samsung-Galaxy-TabPRO-12.2-LTE.html' · '/Samsung-Galaxy-S4-Active.html' ·
'/Samsung-Galaxy-Tab-Active-3.html' · '/Samsung-Galaxy-Tab-S3-9.7.html' · '/Samsung-Galaxy-S6-edge.html' ·
'/Samsung-Galaxy-Note-4-Exynos.html' · '/Samsung-Galaxy-Round.html' ·
'/Samsung-Galaxy-Note-20-Ultra-5G.html' · '/Samsung-ATIV-Q.html' · '/Samsung-ATIV-Smart-PC-PRO.html' ·
'/Samsung-Galaxy-S22-Ultra12-128GB.html' · '/Samsung-Galaxy-Z-Flip-5G.html' · '/Samsung-Galaxy-Z-Flip.html' ·
'/Samsung-Galaxy-Tab-S8-Ultra-5G.html' · '/Samsung-Galaxy-S21-Ultra-16-512GB.html' ·
'/Samsung-Galaxy-S10-Plus-Ram-12GB.html' · '/Samsung-Galaxy-Z-Fold-3.html' · '/Samsung-Galaxy-Z-Fold4.html' ·
'/Samsung-Galaxy-Z-Fold-2-5G.html'

```

```
full_links <- paste0("https://specphone.com",links)
```

## full\_links

'https://specphone.com/Samsung-Galaxy-M13.html' · 'https://specphone.com/Samsung-Galaxy-A23.html' ·  
'https://specphone.com/Samsung-Galaxy-A13.html' · 'https://specphone.com/Samsung-Galaxy-M32-5G.html' ·  
'https://specphone.com/Samsung-Galaxy-A12-Nacho.html' ·  
'https://specphone.com/Samsung-Galaxy-Pocket-Neo.html' ·  
'https://specphone.com/Samsung-Galaxy-Young.html' · 'https://specphone.com/Samsung-Galaxy-J1-Mini.html' ·  
'https://specphone.com/Samsung-Galaxy-A01-Core-1-16GB.html' ·  
'https://specphone.com/Samsung-Galaxy-V-PLUS.html' · 'https://specphone.com/Samsung-Galaxy-Young-2.html' ·  
'https://specphone.com/Samsung-Galaxy-M02.html' · 'https://specphone.com/Samsung-Galaxy-A11.html' ·  
'https://specphone.com/Samsung-Galaxy-J2-Pro-2018.html' ·  
'https://specphone.com/Samsung-Galaxy-A12-2021.html' ·  
'https://specphone.com/Samsung-Galaxy-A21s-3-32GB.html' · 'https://specphone.com/Samsung-Galaxy-J5.html' ·  
'https://specphone.com/Samsung-Galaxy-J4.html' · 'https://specphone.com/Samsung-Galaxy-Core-2-Duos.html' ·  
'https://specphone.com/Samsung-Galaxy-Ace-Plus.html' · 'https://specphone.com/Samsung-Galaxy-A20.html' ·  
'https://specphone.com/Samsung-Galaxy-Chat.html' · 'https://specphone.com/Samsung-Galaxy-Gio.html' ·  
'https://specphone.com/Samsung-Galaxy-Tab-A7-Lite-LTE.html' ·  
'https://specphone.com/Samsung-Galaxy-Tab-A-10.5WIFI.html' ·  
'https://specphone.com/Samsung-Galaxy-Alpha.html' · 'https://specphone.com/Samsung-Galaxy-S3-Slim.html' ·  
'https://specphone.com/Samsung-Galaxy-S4-zoom.html' ·  
'https://specphone.com/Samsung-Galaxy-Xcover-2.html' ·  
'https://specphone.com/Samsung-Galaxy-Tab-8.9-3G-16GB.html' ·  
'https://specphone.com/Samsung-Galaxy-Tab-A8-LTE-2021.html' ·  
'https://specphone.com/Samsung-Galaxy-A8-2018.html' ·  
'https://specphone.com/Samsung-Galaxy-Tab4-8.0-wifi.html' ·  
'https://specphone.com/Samsung-Galaxy-M33-5G.html' · 'https://specphone.com/Samsung-Galaxy-A50.html' ·  
'https://specphone.com/Samsung-Galaxy-E7.html' · 'https://specphone.com/Samsung-Galaxy-S6.html' ·  
'https://specphone.com/Samsung-Galaxy-S20-FE.html' ·  
'https://specphone.com/Samsung-Galaxy-Tab-S4-WIFI.html' · 'https://specphone.com/Samsung-Galaxy-S7.html' ·  
'https://specphone.com/Samsung-Galaxy-Note-5-Exynos.html' ·  
'https://specphone.com/Samsung-Galaxy-TabPRO-12.2-LTE.html' ·  
'https://specphone.com/Samsung-Galaxy-S4-Active.html' ·  
'https://specphone.com/Samsung-Galaxy-Tab-Active-3.html' ·  
'https://specphone.com/Samsung-Galaxy-Tab-S3-9.7.html' ·  
'https://specphone.com/Samsung-Galaxy-S6-edge.html' ·  
'https://specphone.com/Samsung-Galaxy-Note-4-Exynos.html' ·  
'https://specphone.com/Samsung-Galaxy-Round.html' ·  
'https://specphone.com/Samsung-Galaxy-Note-20-Ultra-5G.html' · 'https://specphone.com/Samsung-ATIV-Q.html' ·  
'https://specphone.com/Samsung-ATIV-Smart-PC-PRO.html' ·  
'https://specphone.com/Samsung-Galaxy-S22-Ultra12-128GB.html' ·  
'https://specphone.com/Samsung-Galaxy-Z-Flip-5G.html' · 'https://specphone.com/Samsung-Galaxy-Z-Flip.html' ·  
'https://specphone.com/Samsung-Galaxy-Tab-S8-Ultra-5G.html' ·  
'https://specphone.com/Samsung-Galaxy-S21-Ultra-16-512GB.html' ·  
'https://specphone.com/Samsung-Galaxy-S10-Plus-Ram-12GB.html' ·  
'https://specphone.com/Samsung-Galaxy-Z-Fold-3.html' · 'https://specphone.com/Samsung-Galaxy-Z-Fold4.html' ·  
'https://specphone.com/Samsung-Galaxy-Z-Fold-2-5G.html'



```

result <- data.frame()

for (links in full_links[1:5]) {
  ss_topic <- links %>%
    read_html() %>%
    html_nodes("div.topic") %>%
    html_text2()

  ss_detail <- links %>%
    read_html() %>%
    html_nodes("div.detail") %>%
    html_text2()

  tmp <- data.frame(attribute = ss_topic,
                    value = ss_detail)

  result <- bind_rows(result, tmp)
  print("Progress...")
}

print(result)

```

```

[1] "Progress..."
[1] "Progress..."
[1] "Progress..."
[1] "Progress..."
[1] "Progress..."
      attribute
1   วันเปิดตัว
2   วันวางจำหน่าย
3   ขนาด
4   น้ำหนัก
5   วัสดุ
6   SIM
7   Technology
8   2G
9   3G
10  4G
11  5G
12  ความเร็ว
13  ประเภท
14  ความจุหน่วยความจำ

```

```
print(head(result),3)
```

attribute

value

1	วันเปิดตัว	มิถุนายน 2565
2	วันวางจำหน่าย	ยังไม่วางจำหน่าย
3	ขนาด	165.40 x 76.90 x 8.40 มม.
4	น้ำหนัก	192 กรัม
5	วัสดุ	Glass front, plastic back, plastic frame
6	SIM	รองรับ 2 ซิมการ์ด (nano sim, nano sim)