# Mini Project 01 - IMDB web *scraping*

```r
library(tidyverse)
library(rvest) #scratch data from the internet
```

Attaching package: 'rvest'

The following object is masked from 'package:readr':

    guess_encoding

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

```r
print(url)
```

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

```r
#readhtml
imdb <- read_html(url)
```

```r
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 .
```

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

```
#movie title
ratings <- imdb %>%
 html_nodes("div.ratings-imdb-rating") %>%
 html_text2() %>%
 as.numeric()
```

```
ratings[1:10]
```

9.3 · 9.2 · 9 · 9 · 9 · 9 · 9 · 8.9 · 8.8 · 8.8

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

```
num_votes
```

'Votes: 2,693,841 | Gross: $28.34M | Top 250: #1' · 'Votes: 1,869,358 | Gross: $134.97M | Top 250: #2' ·
'Votes: 1,362,436 | Gross: $96.90M | Top 250: #6' · 'Votes: 2,667,687 | Gross: $534.86M | Top 250: #3' ·
'Votes: 795,929 | Gross: $4.36M | Top 250: #5' · 'Votes: 1,278,157 | Gross: $57.30M | Top 250: #4' ·
'Votes: 1,855,813 | Gross: $377.85M | Top 250: #7' · 'Votes: 2,067,755 | Gross: $107.93M | Top 250: #8' ·
'Votes: 2,366,510 | Gross: $292.58M | Top 250: #14' · 'Votes: 2,139,302 | Gross: $37.03M | Top 250: #12' ·
'Votes: 1,885,310 | Gross: $315.54M | Top 250: #9' · 'Votes: 2,092,010 | Gross: $330.25M | Top 250: #11' ·
'Votes: 766,107 | Gross: $6.10M | Top 250: #10' · 'Votes: 1,675,772 | Gross: $342.55M | Top 250: #13' ·
'Votes: 1,168,879 | Gross: $46.84M | Top 250: #17' · 'Votes: 1,923,105 | Gross: $171.48M | Top 250: #16' ·
'Votes: 1,012,860 | Gross: $112.00M | Top 250: #18' · 'Votes: 1,299,410 | Gross: $290.48M | Top 250: #15' ·
'Votes: 1,849,404 | Gross: $188.02M | Top 250: #25' · 'Votes: 1,663,075 | Gross: $100.13M | Top 250: #19' ·
'Votes: 1,440,814 | Gross: $130.74M | Top 250: #22' · 'Votes: 1,371,850 | Gross: $322.74M | Top 250: #28' ·
'Votes: 1,309,748 | Gross: $136.80M | Top 250: #27' · 'Votes: 1,399,262 | Gross: $216.54M | Top 250: #24' ·
'Votes: 699,802 | Gross: $57.60M | Top 250: #26' · 'Votes: 1,105,607 | Gross: $204.84M | Top 250: #29' ·
'Votes: 770,215 | Gross: $10.06M | Top 250: #31' · 'Votes: 760,963 | Gross: $7.56M | Top 250: #23' ·
'Votes: 466,075 | Top 250: #21' · 'Votes: 348,476 | Gross: $0.27M | Top 250: #20' · 'Votes: 58,591 | Top 250: #45' ·
'Votes: 872,895 | Gross: $13.09M | Top 250: #42' · 'Votes: 815,793 | Gross: $53.37M | Top 250: #34' ·
'Votes: 1,213,822 | Gross: $210.61M | Top 250: #30' · 'Votes: 1,509,148 | Gross: $187.71M | Top 250: #37' ·
'Votes: 1,341,245 | Gross: $53.09M | Top 250: #41' · 'Votes: 1,332,787 | Gross: $132.38M | Top 250: #39' ·
'Votes: 672,241 | Gross: $83.47M | Top 250: #53' · 'Votes: 1,168,761 | Gross: $19.50M | Top 250: #35' ·
'Votes: 888,783 | Gross: $78.90M | Top 250: #51' · 'Votes: 1,090,842 | Gross: $23.34M | Top 250: #40' ·
'Votes: 1,128,403 | Gross: $6.72M | Top 250: #38' · 'Votes: 838,403 | Gross: $32.57M | Top 250: #32' ·
'Votes: 1,065,127 | Gross: $422.78M | Top 250: #36' · 'Votes: 865,026 | Gross: $13.18M | Top 250: #46' ·
'Votes: 575,536 | Gross: $1.02M | Top 250: #43' · 'Votes: 332,645 | Gross: $5.32M | Top 250: #48' ·
'Votes: 676,423 | Gross: $32.00M | Top 250: #33' · 'Votes: 280,794 | Top 250: #44' ·
'Votes: 263,911 | Gross: $11.99M | Top 250: #50'

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

head(df)
```

A data.frame: 6 × 3

| | title | rating | num_vote |
|---|---|---|---|
| | <chr> | <dbl> | <chr> |
| 1 | 1. The Shawshank Redemption (1994) | 9.3 | Votes: 2,693,841 | Gross: $28.34M | Top 250: #1 |
| 2 | 2. The Godfather (1972) | 9.2 | Votes: 1,869,358 | Gross: $134.97M | Top 250: #2 |
| 3 | 3. Schindler's List (1993) | 9.0 | Votes: 1,362,436 | Gross: $96.90M | Top 250: #6 |
| 4 | 4. The Dark Knight (2008) | 9.0 | Votes: 2,667,687 | Gross: $534.86M | Top 250: #3 |
| 5 | 5. 12 Angry Men (1957) | 9.0 | Votes: 795,929 | Gross: $4.36M | Top 250: #5 |
| 6 | 6. The Godfather Part II (1974) | 9.0 | Votes: 1,278,157 | Gross: $57.30M | Top 250: #4 |

# Mini Project 0. - Specphone Phone database

```r
library(tidyverse)
library(rvest) #scratch data from the internet
```

```
Warning message in system("timedatectl", intern = TRUE):
"running command 'timedatectl' had status 1"
Warning message:
"Failed to locate timezone database"
── Attaching packages ──────────────────────────── tidyverse 1.3.1

✓ ggplot2 3.3.5      ✓ purrr   0.3.4
✓ tibble  3.1.5      ✓ dplyr   1.0.7
✓ tidyr   1.1.4      ✓ stringr 1.4.0
✓ readr   2.0.2      ✓ forcats 0.5.1

── Conflicts ──────────────────────────────── tidyverse_conflicts()
✗ dplyr::filter()  masks stats::filter()
✗ purrr::flatten() masks jsonlite::flatten()
✗ dplyr::lag()     masks stats::lag()


Attaching package: 'rvest'
```

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

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

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

```
data.frame(attribute = att, value = value)
```

A data.frame: 31 × 2

| attribute | value |
|---|---|
| <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) |
| ความละเอียดวีดีโอ | 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 sumsung

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")
```

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

```
result <- data.frame()

for (link in full_links[1:10]){
    ss_topic <- link %>%
    read_html() %>%
 html_nodes("div.topic") %>%
 html_text2()

    ss_detail<- link %>%
    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...."
[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
```

```
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)
```

```
#write csv
write_csv(result,"result_ss_phone.csv")
```