

# Mini Project 01 - 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"
```

```
# read html
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" width .
```

```
# movie title
titles <- imdb %>%
  html_nodes("h3.lister-item-header") %>%
  html_text2() # text2 จะลบอักขระพิเศษด้วย
```

```
# rating
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

```
# number of votes
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
)

head(df)
```

A data.frame: 6 × 3

	title	rating	num_vote
	<chr>	<dbl>	<chr>
1	1. The Shawshank Redemption (1994)	9.3	Votes: 2,696,378   Gross: \$28.34M   Top 250: #1
2	2. The Godfather (1972)	9.2	Votes: 1,871,055   Gross: \$134.97M   Top 250: #2
3	3. Schindler's List (1993)	9.0	Votes: 1,363,253   Gross: \$96.90M   Top 250: #6
4	4. The Dark Knight (2008)	9.0	Votes: 2,670,071   Gross: \$534.86M   Top 250: #3
5	5. 12 Angry Men (1957)	9.0	Votes: 796,245   Gross: \$4.36M   Top 250: #5
6	6. The Godfather Part II (1974)	9.0	Votes: 1,278,949   Gross: \$57.30M   Top 250: #4

# Mini Project 02 - Specphone Phone Database

```
library(tidyverse)
library(rvest)
```

```
url <- read_html("https://specphone.com/Xiaomi-Redmi-10A.html")
```

```
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	มกราคม 2566
วันวางจำหน่าย	ยังไม่วางจำหน่าย
ขนาด	164.90 x 77.10 x 9.00 มม.
น้ำหนัก	194 กรัม
วัสดุ	ไมโครรับ
SIM	รองรับ 2 ซิมการ์ด (nano sim, nano sim)
Technology	HSPA, LTE
2G	850/900/1800/1900
3G	850/900/1900/2100
4G	850/900/1900/2100/2600
5G	-
ความเร็ว	HSPA, LTE
ประเภท	IPS LCD
ขนาดหน้าจอ	6.53 นิ้ว
ความละเอียด	720 x 1600 pixels
ระบบปฏิบัติการ	Android 11
ชิปประมวลผล	MediaTek Helio G25 2 GHz
ชิปกราฟิก	PowerVR GE8320
หน่วยความจำ	3 GB
ความจุ	32 GB
Memory Card	ไมโครรับ
กล้องหลัก	ตัวที่ 1: 13 MP, f/2.2, (wide), 1.0µm, AF
ความละเอียดวิดีโอ	1080p@30fps
กล้องหน้า	ตัวที่ 1: 5 MP
Bluetooth	5.0, A2DP, LE
Wi-Fi	802.11 a/b/g/n/ac, dual-b
USB	micro USB
GPS	GPS, GLONASS, GALILEO, BD
NFC	ไมโครรับ
ความจุ	5,000 mAh
ประเภท	Non-removable Li-Po Batt

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

```
# link to all samsung smart phone
links <- samsung_url %>%
  html_nodes("li.mobile-brand-item a") %>% # เทคนิค " a" จงหาตัว a ที่อยู่ใน li.mobile
  html_attr("href")
```

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

```
full_links
```

```
result <- data.frame()

for (link in full_links[1:5]) {
  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(result)
```

```
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")
```