

Mini Project 01 -IMDB web scraping

```
library(tidyverse)
library(rvest) #scrap data from internet
library(dplyr)
```

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

```
— Attaching packages — tidyverse 1.3.2
```

```
✓ 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'
```

```
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/fb
[1] <head>\n<meta http-equiv="Content-Type" content="text/html; charset=UTF-
[2] <body id="styleguide-v2" class="fixed">\n          <img height="1" wid
```

```
##movie title
titles <- imdb %>%
  html_nodes("h3.lister-item-header") %>%
  html_text2() ##can use text and text 2
```

```
titles[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)'
```

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

```
ratings
```

```
9.3 · 9.2 · 9 · 9 · 9 · 9 · 9 · 9 · 8.9 · 8.8 · 8.8 · 8.8 · 8.8 · 8.8 · 8.8 · 8.7 · 8.7 · 8.7 · 8.7 · 8.6 · 8.6 · 8.6 · 8.6 · 8.6 · 8.6 · 8.6 · 8.6 ·
8.6 · 8.6 · 8.6 · 8.6 · 8.6 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5 · 8.5
```

```
##number of vote
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: 2,675,547 Gross: \$28.34M Top 250: #1
2	2. The Godfather (1972)	9.2	Votes: 1,854,564 Gross: \$134.97M Top 250: #2
3	3. The Dark Knight (2008)	9.0	Votes: 2,648,727 Gross: \$534.86M Top 250: #3
4	4. The Lord of the Rings: The Return of the King (2003)	9.0	Votes: 1,843,787 Gross: \$377.85M Top 250: #7
5	5. Schindler's List (1993)	9.0	Votes: 1,354,286 Gross: \$96.90M Top 250: #6
6	6. The Godfather Part II (1974)	9.0	Votes: 1,269,654 Gross: \$57.30M Top 250: #4

Mini project 02 - Specphone Phone Database

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

```
url <- "https://specphone.com/Huawei-Mate-50.html"
```

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

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

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

attribute	value
<chr>	<chr>
วันเปิดตัว	พฤศจิกายน 2565
วันวางจำหน่าย	พฤศจิกายน 2565, ยังไม่วางจำหน่าย
ขนาด	161.50 x 75.10 x 8.00 มม.
น้ำหนัก	202 กรัม
วัสดุ	Glass front, glass back or eco leather back, aluminum frame
SIM	รองรับ 2 ซิมการ์ด (nano sim, nano sim)
Technology	HSPA, LTE-A
2G	850/900/1800/1900
3G	850/900/1900/2100
4G	850/900/1900/2100/2600
5G	-
ความเร็ว	HSPA, LTE-A
ประเภท	OLED
ขนาดหน้าจอ	6.70 นิ้ว
ความละเอียด	1224 x 2700 pixels
ระบบปฏิบัติการ	EMUI 13
ชิปประมวลผล	Qualcomm Snapdragon 8+ Gen 1 SM8475 3.19 GHz
ชิปกราฟิก	Adreno 730
หน่วยความจำ	8 GB
ความจุ	256 GB
Memory Card	microSD (256)
กล้องหลัก	ตัวที่ 1: 50 MP, f/1.4-f/4.0, 24mm (wide), PDAF, Laser AF, OIS ตัวที่ 2: 12 MP, f/3.4, 125mm (periscope telephoto), PDAF, OIS, 5x optical zoom ตัวที่ 3: 13 MP, f/2.2, 13mm, 120° (ultrawide), PDAF
ความละเอียดวิดีโอ	4K@30/60fps, 1080p@30/60/120/240fps, 1080p@960fps, gyro-EIS
กล้องหน้า	ตัวที่ 1: 13 MP, f/2.4, 18mm (ultrawide)
Bluetooth	5.2, A2DP, LE
Wi-Fi	802.11 a/b/g/n/ac/6, dual
USB	Type-C
GPS	with dual-band A-GPS, GLO
อื่นๆ	~

```
##all samsung 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") %>% ##หาตัว A ที่อยู่ใน li
  html_attr("href") #หา attribute 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' ·

```
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(attributes=ss_topic,
                    value= ss_detail)
  result <- bind_rows(result,tmp)
```

```
print("progress...")

}

print(result)
```

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

	attributes	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(result,"result_ss_phone.csv")
```