

Web scraping in R

Ujang Fahmi

Tujuan

1. Apa web Scraping?
2. Mengapa kita perlu melakukan web Scraping?
3. Bagaimana cara melakukan web Scraping?
4. Prinsip Kerja web Scraping di R dengan paket rvest

-
1. Scraping table
 2. Scraping static webpage
 3. Scraping multipage
 4. Storing scraped data in data frame

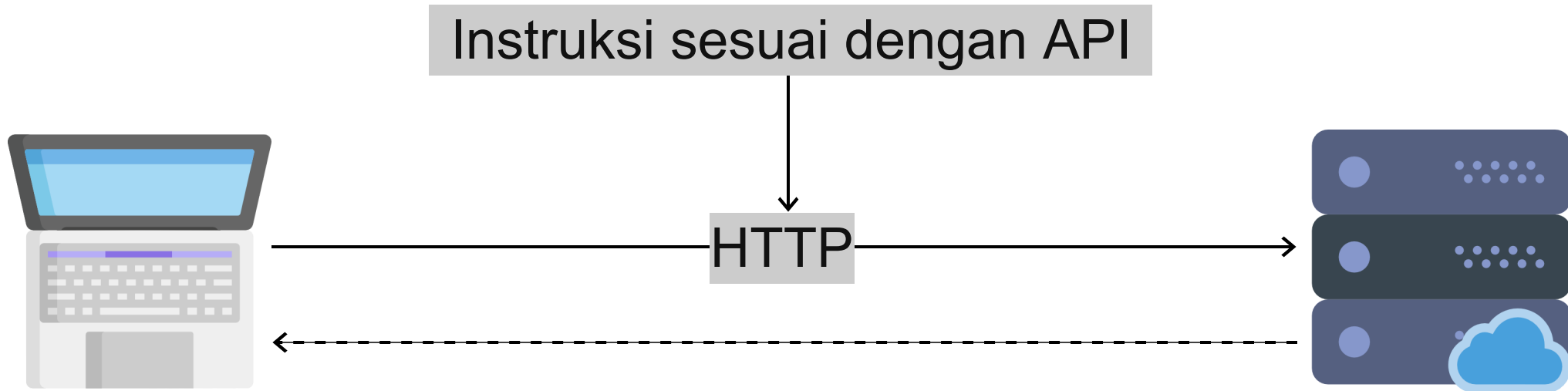


UMUM

KHUSUS

Apa web Scraping?

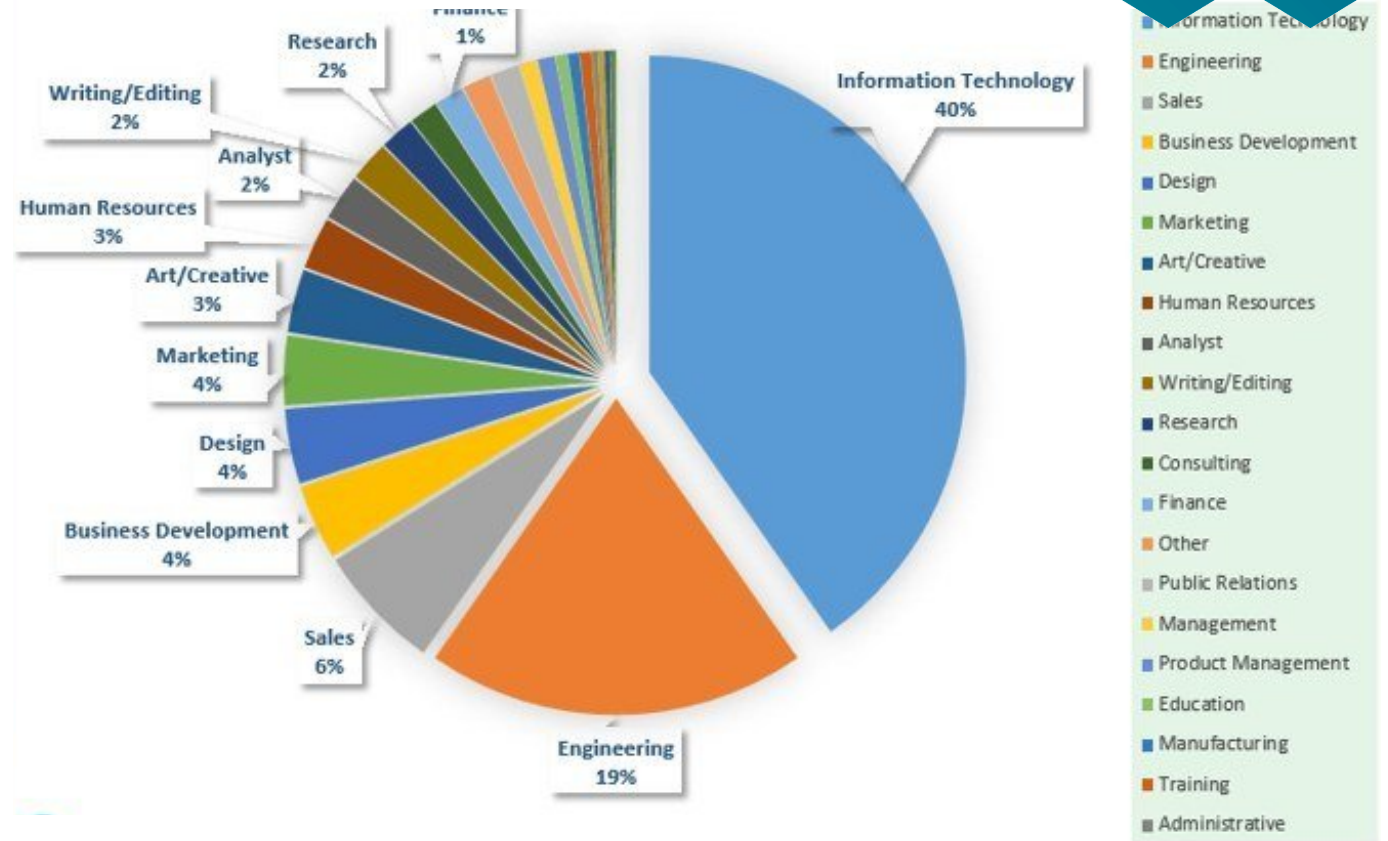
Web scraping, data scraping, web *extraction* adalah proses mengekstrak data dari halaman sebuah website dengan memanfaatkan beberapa teknik seperti *copy paste*, *html parsing*, *DOM parsing*, dan lain-lain.



Mengapa Scraping diperlukan?

Data is new differentiator
Bernard Marr, Forbes

An essential skill to
acquire in today's digital
world



Bagaimana Scraping di R?

1. Download html using `read_html()`
2. Extract specific nodes using `html_nodes()`
3. Extract element of a nodes using `html_text()`, `html_attr()`, `html_table()`, etc.
4. Pre-process



Menggunakan **xpath** dan css nodes



```
1 <?xml version="1.0" encoding="UTF-8"?>
2
3 <bookstore>
4
5 <book>
6   <title lang="en">Harry Potter</title>
7   <price>29.99</price>
8 </book>
9
10 <book>
11   <title lang="en">Learning XML</title>
12   <price>39.95</price>
13 </book>
14
15 </bookstore>
```

Expression	Result
<i>nodename</i>	Selects all nodes with the name " <i>nodename</i> "
/	Selects from the root node
//	Selects nodes in the document from the current node that match the selection no matter where they are
.	Selects the current node
..	Selects the parent of the current node
@	Selects attributes

Menggunakan **xpath** dan css nodes



Path Expression	Result
bookstore	Selects all nodes with the name "bookstore"
/bookstore	Selects the root element bookstore Note: If the path starts with a slash (/) it always represents an absolute path to an element!
bookstore/book	Selects all book elements that are children of bookstore
//book	Selects all book elements no matter where they are in the document
bookstore//book	Selects all book elements that are descendant of the bookstore element, no matter where they are under the bookstore element
//@lang	Selects all attributes that are named lang
//book/title //book/price	Selects all the title AND price elements of all book elements

Contoh `xpath`

Mendapatkan judul artikel bagian *most shared* dari laman <https://www.kdnuggets.com/news/top-stories.html>



```
1 library(tidyverse)
2 library(rvest)
3
4 page <- read_html("https://www.kdnuggets.com/news/top-stories.html")
5
6 # xpath
7 page %>%
8   html_nodes(xpath = "//ol[1]/li/a/b") %>%
9   html_text(trim = TRUE)
10
11 page %>%
12   html_nodes(xpath = "//ol[1]//li//a") %>%
13   html_attr("href") %>%
14   ifelse(. == " ", NA, .)
```

Menggunakan xpath dan **css**nodes



Selector	Example	Example description
<i>.class</i>	.intro	Selects all elements with class="intro"
<i>.class1.class2</i>	<div class="name1 name2">...</div>	Selects all elements with both <i>name1</i> and <i>name2</i> set within its class attribute
<i>element>element</i>	div > p	Selects all <p> elements where the parent is a <div> element
<i>element+element</i>	div + p	Selects all <p> elements that are placed immediately after <div> elements
<i>:nth-child(n)</i>	p:nth-child(2)	Selects every <p> element that is the second child of its parent
<i>element1~element2</i>	p ~ ul	Selects every element that are preceded by a <p> element
<i>element,element</i>	div, p	Selects all <div> elements and all <p> elements

Contoh

Menggunakan css selector



```
1 library(tidyverse)
2 library(rvest)
3
4 page <- read_html("https://www.kdnuggets.com/news/top-stories.html")
5
6 page %>%
7   html_nodes(css = "ol:nth-child(3) > li > a > b") %>%
8   html_text(trim = TRUE) %>%
9   data_frame()
10
11 page %>%
12   html_nodes(css = "ol:nth-child(3) > li > a") %>%
13   html_attr("href") %>%
14   data_frame()
```

Multiple Page 1



halaman → <https://jdih.dprd-diy.go.id/?cat=5>

halaman 1 → <https://jdih.dprd-diy.go.id/?pagenum=2&totalrow=27&cat=5>

halaman 2 → <https://jdih.dprd-diy.go.id/?pagenum=2&totalrow=27&cat=5>

baseurl

<https://jdih.dprd-diy.go.id/?pagenum=>

nomor halaman

angka halaman

endurl

[2&totalrow=27&cat=5](https://jdih.dprd-diy.go.id/?pagenum=2&totalrow=27&cat=5)

Multiple Page 2



```
1 library(rvest)
2
3 baseurls <- "https://jdih.dprd-diy.go.id/?pagenum="
4 page <- (0:2)
5 endurls <- "&totalrow=27&cat=5"
6
7 # list kosong
8 urls <- list()
9
10 # membuat daftar urls
11 for (i in seq_along(page)) {
12   url<- paste0(baseurls, page[i], endurls)
13   urls[[i]] <- url
14 }
```

Target: Membuat daftar url
yang akan diambil datanya



```
1 # list kosong untuk menampung hasil
2 undang2 <- list()
3
4 # loop over the urls and get the table from each page
5 for (i in seq_along(urls)) {
6   pages <- read_html(urls[[i]])
7
8   tentang <- pages %>%
9     html_nodes("tr:nth-child(3) td~ td+ td") %>%
10    html_text(trim = TRUE) %>%
11    data_frame()
12
13   download <- pages %>%
14     html_nodes("tr:nth-child(6) > td") %>%
15     html_node("a") %>%
16     html_attr("href") %>%
17     ifelse(. == " ", NA, .) %>%
18     data_frame()
19
20   download$. <- paste0("https://jdih.dprd-diy.go.id/", download$.)
21   undang <- bind_cols(tentang, download)
22   undang2[[i]] <- undang
23 }
```



Target: mengambil elemen
teks tentang dan download
dari tiap halaman

Scraping tabel



```
1 library(rvest)
2 library(tidyverse)
3
4 pages <- read_html("http://www.dpr.go.id/jdih/pp")
5
6 hasil <- html_nodes(pages, "table") [[1]] %>%
7   html_table()
8
9 hasil <- pages %>%
10   html_nodes("table") %>%
11   html_table()
12
13 hasil <- hasil[[1]]
14
15 pages %>%
16   html_nodes("table") %>%
17   html_table()
```

Target: mengambil tabel dan isinya dari sebuah halaman website

Rangkuman



1. Scraping bisa dilakukan di R
2. Scraping memanfaatkan css nodes dan atau xpath nodes
3. Tidak ada nodes yang lebih baik, terkadang kita harus mencobanya dan melihat hasilnya
4. Melihat perubahan urls merupakan strategi umum yang bisa digunakan untuk melakukan scraping dari beberapa halaman
5. Nodes untuk tabel umumnya adalah table, jika banyak bisa diurukan tabel ke berapa
6. Sering mencoba (praktik)