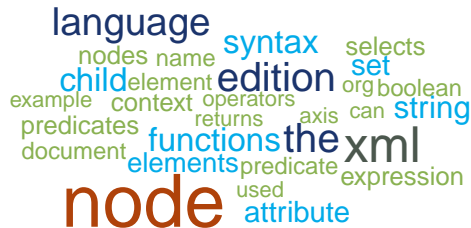# Introduction to Web Scraping with R

Scraping Multiple Pages



Simon Munzert | IPSDS

# An advanced scraping scenario

# Motivation

- until now, the toy examples were limited to single HTML pages

- often, we want to scrape data from multiple pages

- in such scenarios, automating the scraping process becomes **really** powerful

- also, the principles of polite scraping are more relevant

# The scenario

**Goal:** examine download statistics of articles of the Journal of Statistical Software

- download HTML pages
- extract bibliometrical information

**Tasks:**

- identify relevant resources on http://www.jstatsoft.org/
- download HTML pages
- import them into R
- extract information via XPath

# Scraping multiple pages with R

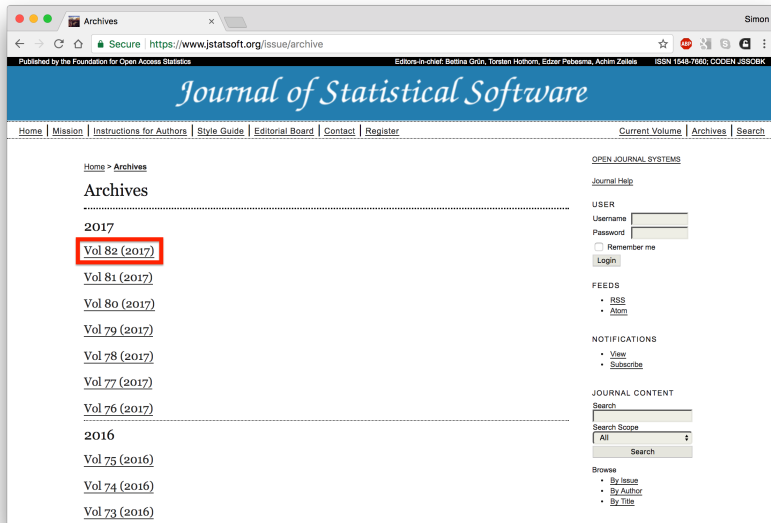# Step 1: Inspect the source



## Procedure

- source: http://www.jstatsoft.org/
- go to "Archives"

# Step 1: Inspect the source



## Procedure

- source: http://www.jstatsoft.org/
- go to "Archives"
- inspect the most recent volume

# Step 1: Inspect the source



## Procedure

- source: http://www.jstatsoft.org/
- go to "Archives"
- inspect the most recent volume
- inspect the first article

# Step 1: Inspect the source

## Procedure

- source: http://www.jstatsoft.org/
- go to "Archives"
- inspect the most recent volume
- inspect the first article

# Step 1: Inspect the source



## Procedure

- source: http://www.jstatsoft.org/
- go to "Archives"
- inspect the most recent volume
- inspect the first article
- inspect the page views element

# Step 1: Inspect the source

## Procedure

- source: http://www.jstatsoft.org/

- go to "Archives"

- inspect the most recent volume

- inspect the first article

- inspect the page views element

- it's in a table!

# Step 2: Develop a scraping strategy

## Observations

- getting the information out of the table will be straightforward
- this applies to all articles (check other articles on a sample basis)
- what we need is the set of **URLs leading to all articles**

# Step 2: Develop a scraping strategy

## Observations

- getting the information out of the table will be straightforward
- this applies to all articles (check other articles on a sample basis)
- what we need is the set of **URLs leading to all articles**

## Inspecting the URLs

- the URL of the selected article looks as follows:
  https://www.jstatsoft.org/article/view/v082i01
- we find out that the final part, v082i01, always follows the same pattern:
  v<volume number>i<issue number>

# Step 2: Develop a scraping strategy

## Let's try to construct the list of URLs from scratch

```
R code
1   baseurl <- "http://www.jstatsoft.org/article/view/v"
2   volurl <- paste0("0", seq(1, 78, 1))
3   volurl[1:9] <- paste0("00", seq(1, 9, 1))
4   brurl <- paste0("0", seq(1, 9, 1))
5   urls_list <- paste0(baseurl, volurl)
6   urls_list <- paste0(rep(urls_list, each = 9), "i", brurl)
7   urls_list[1:5]
    [1] "http://www.jstatsoft.org/article/view/v001i01"
    [2] "http://www.jstatsoft.org/article/view/v001i02"
    [3] "http://www.jstatsoft.org/article/view/v001i03"
    [4] "http://www.jstatsoft.org/article/view/v001i04"
    [5] "http://www.jstatsoft.org/article/view/v001i05"
8   names <- paste0(rep(volurl, each = 9), "_", brurl, ".html")
9   names[1:5]
    [1] "001_01.html" "001_02.html" "001_03.html" "001_04.html" "001_05.html"
                                                                              end
```

# Step 3: Download the files

## Set working directory

R code

```
10   tempwd <- ("data/jstatsoftStats")
11   dir.create(tempwd)
12   setwd(tempwd)
```
end

## Download pages

R code

```
13   folder <- "html_articles/"
14   dir.create(folder)
15   for (i in 1:length(urls_list)) {
16       if (!file.exists(paste0(folder, names[i]))) {
17           download.file(urls_list[i], destfile = paste0(folder, names[i]))
18           Sys.sleep(runif(1, 0, 1))
19       }
20   }
```
end

# Step 3: Download the files

## Check success

R code

```
21  list_files <- list.files(folder, pattern = "0.*")
22  list_files_path <- list.files(folder, pattern = "0.*", full.names = TRUE)
23  length(list_files)

    [1] 666
```

end

# Step 4: Import files and parse out information

## Build loop

R code

```
24  authors <- character()
25  title <- character()
26  statistics <- character()
27  numViews <- numeric()
28  datePublish <- character()
29  for (i in 1:length(list_files_path)) {
30      html_out <- read_html(list_files_path[i])
31      table_out <- html_table(html_out, fill = TRUE)[[6]]
32      authors[i] <- table_out[1, 2]
33      title[i] <- table_out[2, 2]
34      statistics[i] <- table_out[4, 2]
35      numViews[i] <- statistics[i] %>% str_extract("[[:digit:]]+") %>% as.numeric()
36      datePublish[i] <- statistics[i] %>% str_extract("[[:digit:]]{4}-[[:digit:]]{2}-[[:digit
    :]]{2}.$") %>%
37          str_replace("\\.", "")
38  }
```

end

# Step 4: Import files and parse out information

## Inspect parsed data

```
R code
39  authors[1:3]
    [1] "Ronald Barry"   "Jason Bond, George Michailides"   "Thomas Lumley"
40  title[1:2]
    [1] "A Diagnostic to Assess the Fit of a Variogram Model to Spatial Data"
    [2] "Homogeneity Analysis in Xlisp-Stat"
41  numViews[1:3]
    [1] 5835 3939 4379
                                                                              end
```

## Construct data frame

```
R code
42  dat <- data.frame(authors = authors, title = title, numViews = numViews, datePublish =
    datePublish)
43  dim(dat)
    [1] 666    4
                                                                              end
```
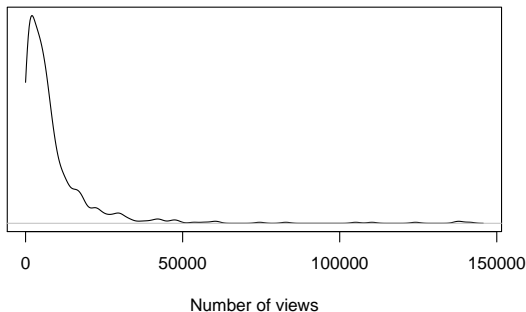
# Step 5: Visualize data

## Density plot of download statistics

R code

```
44  plot(density(dat$numViews, from = 0), yaxt="n", ylab="", xlab="Number of views", main="
    Distribution of article page views in JStatSoft")
```
end

**Distribution of article page views in JStatSoft**



Number of views

# Summary

# Summary

- scraping data from multiple pages is no problem in R

- most of the brain work often goes into developing a scraping strategy and tidying the data, not into the actual downloading/scraping part

- scraping is also possible in even more complex scenarios, e.g., when HTML forms are involved or you have to take care of cookies or authentication

- this is beyond the scope of this course → check out the book for more applications



Source: Horia Varlan