

Lab 06 - Regular Expressions and Web Scraping

Learning goals

- Use a real world API to make queries and process the data.
- Use regular expressions to parse the information.
- Practice your GitHub skills.

Lab description

In this lab, we will be working with the NCBI API to make queries and extract information using XML and regular expressions. For this lab, we will be using the `httr`, `xml2`, and `stringr` R packages.

This markdown document should be rendered using `github_document` document ONLY and pushed to your *JSC370-labs* repository in `lab06/README.md`.

Question 1: How many sars-cov-2 papers?

Build an automatic counter of sars-cov-2 papers using PubMed. You will need to apply XPath as we did during the lecture to extract the number of results returned by PubMed in the following web address:

<https://pubmed.ncbi.nlm.nih.gov/?term=sars-cov-2>

Complete the lines of code:

```
url <- "https://pubmed.ncbi.nlm.nih.gov/?term=sars-cov-2"
# Downloading the website
website <- xml2::read_html(url)
# Finding the counts
counts <- xml2::xml_find_first(website, "/html/body/main/div[9]/div[2]/div[2]/div[1]/div[1]/span")
# Turning it into text
counts <- as.character(counts)
# Extracting the data using regex
stringr::str_extract(counts, "[\\d,]+")
```

- How many sars-cov-2 papers are there?

There are 192,677 sars-cov-2 papers

Don't forget to commit your work!

Question 2: Academic publications on COVID19 and Hawaii

Use the function `httr::GET()` to make the following query:

1. Baseline URL: `https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi`
2. Query parameters:
 - db: pubmed
 - term: covid19 hawaii
 - retmax: 1000

The parameters passed to the query are documented here.

```
library(httr)
baseurl <- "https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi"
query_ids <- GET(
  url = baseurl,
  query = list(db = "pubmed",
               term = "covid19 hawaii",
               retmax = 1000)
)
# Extracting the content of the response of GET
ids <- httr::content(query_ids)
```

The query will return an XML object, we can turn it into a character list to analyze the text directly with `as.character()`. Another way of processing the data could be using lists with the function `xml2::as_list()`. We will skip the latter for now.

Take a look at the data, and continue with the next question (don't forget to commit and push your results to your GitHub repo!).

Question 3: Get details about the articles

The Ids are wrapped around text in the following way: `<Id>... id number ...</Id>`. we can use a regular expression that extract that information. Fill out the following lines of code:

```
# Turn the result into a character vector
ids <- as.character(ids)
# Find all the ids
ids <- stringr::str_extract_all(ids, "<Id>[\\d]+</Id>")[[1]]
# Remove all the leading and trailing <Id> </Id>. Make use of "/"
ids <- stringr::str_remove_all(ids, "<Id>|</Id>")
```

With the ids in hand, we can now try to get the abstracts of the papers. As before, we will need to coerce the contents (results) to a list using:

1. Baseline url: `https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi`
2. Query parameters:
 - db: pubmed
 - id: A character with all the ids separated by comma, e.g., "1232131,546464,13131"

- retmax: 1000
- rettype: abstract

Pro-tip: If you want `GET()` to take some element literal, wrap it around `I()` (as you would do in a formula in R). For example, the text "123,456" is replaced with "123%2C456". If you don't want that behavior, you would need to do the following `I("123,456")`.

```
baseurl2 = "https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi"
publications <- GET(
  url = baseurl2,
  query = list(
    db="pubmed",
    id= paste(ids, collapse = ","),
    retmax=1000,
    rettype="abstract"
  )
)
# Turning the output into character vector
publications <- httr::content(publications)
publications_txt <- as.character(publications)
```

With this in hand, we can now analyze the data. This is also a good time for committing and pushing your work!

Question 4: Distribution of universities, schools, and departments

Using the function `stringr::str_extract_all()` applied on `publications_txt`, capture all the terms of the form:

1. University of ...
2. ... Institute of ...

Write a regular expression that captures all such instances

```
institution <- stringr::str_extract_all(
  publications_txt,
  "University of [[:alnum:]]+ | [[:alnum:]]+ Institute of [[:alnum:]]+"
)
institution <- unlist(institution)
as.data.frame(table(institution))
```

Repeat the exercise and this time focus on schools and departments in the form of

1. School of ...
2. Department of ...

And tabulate the results

```
schools_and_deps <- stringr::str_extract_all(
  publications_txt,
  "School of [[:alnum:]]+| Department of [[:alnum:]]+"
)
as.data.frame(table(schools_and_deps))
```

Question 5: Form a database

We want to build a dataset which includes the title and the abstract of the paper. The title of all records is enclosed by the HTML tag `ArticleTitle`, and the abstract by `Abstract`.

Before applying the functions to extract text directly, it will help to process the XML a bit. We will use the `xml2::xml_children()` function to keep one element per id. This way, if a paper is missing the abstract, or something else, we will be able to properly match PUBMED IDS with their corresponding records.

```
pub_char_list <- xml2::xml_children(publications)
pub_char_list <- sapply(pub_char_list, as.character)
```

Now, extract the abstract and article title for each one of the elements of `pub_char_list`. You can either use `sapply()` as we just did, or simply take advantage of vectorization of `stringr::str_extract`

```
abstracts <- stringr::str_extract(pub_char_list, "<Abstract>(\\n|.)+</Abstract>")
abstracts <- stringr::str_remove_all(abstracts, "<[/[:alnum:]]+>|<[/[:alnum:]]+>")
abstracts <- stringr::str_replace_all(abstracts, "\\s+", " ")
abstracts <- stringr::str_remove_all(abstracts, "\\n")
table(is.na(abstracts))
```

- How many of these don't have an abstract?

53 of them have no abstract

Now, the title

```
titles <- stringr::str_extract(pub_char_list, "<Title>(\\n|.)+</Title>")
titles <- stringr::str_remove_all(titles, "<[/[:alnum:]]+>|<[/[:alnum:]]+>")
table(is.na(titles))
```

- How many of these don't have a title ?

All of them have title

Finally, put everything together into a single `data.frame` and use `knitr::kable` to print the results

```
database <- data.frame(
  ids, titles, abstracts
)
knitr::kable(database)
```

Done! Knit the document, commit, and push.

Final Pro Tip (optional)

You can still share the HTML document on github. You can include a link in your `README.md` file as the following:

```
View [here](https://cdn.jsdelivr.net/gh/:user/:repo@:tag/:file)
```

For example, if we wanted to add a direct link the HTML page of lecture 6, we could do something like the following:

View Week 6 Lecture [\[here\]](#)()