# Harvesting online data

Short course at Santa Clara University, ECON173

Johannes Fritz

5/5/22, 5/12/22

# WHY AM I HERE?

►Economist with a strong "data science tilt"
> Web harvesting
> Machine learning
> Natural language processing

►Data collection is my day job
> Global Trade Alert
Dataset on industrial policy choices around the world since 2009
> Digital Policy Alert
Dataset on G20+Europe's digital policy choices since 2021

# YOU WILL LEARN 3 TOOLS

1. Harvesting the data:
   Headless **web scraping** using R. — today

2. Parsing the data:
   Introduction to **Regular Expressions**. — next week

3. Organizing the data:
   Basics of **relational database design**.

All material at:
https://github.com/johannesfritz/scu-harvesting

# Getting started

► What projects do you want to tackle?

► What experience do you already have?

► What tools do you want to learn?

# The U.S. Congressional Record

My teaching example for this short course

# THE U.S. CONGRESSIONAL HEARINGS

Source

► https://www.govinfo.gov/app/collection/CHRG/

► or search:
«GPO Congressional Hearings»

# WEBSITE CONTENTS

**Download**

- PDF
- Text
- MODS
- PREMIS
- ZIP

**Actions**

- Browse Congressional Hearings
- Help
- CGP Record

## Content Details

House Hearing, 117th Congress - CLEARING THE AIR: SCIENCE-BASED STRATEGIES TO PROTECT WORKERS FROM COVID-19 INFECTIONS

| Summary | Document in Context |
| --- | --- |

| | |
| --- | --- |
| Category | Congressional Committee Materials |
| Collection | Congressional Hearings |
| SuDoc Class Number | Y 4.ED 8/1:117-1 |
| ILS System ID | 001171899 |
| Congress | 117th |
| Congress Chamber | House of Representatives |
| Hearing Sub Type | General |
| Held Date | March 11, 2021 |
| Committee and Subcommittee | Committee on Education and Labor |
| Members | Joe Wilson (SC); Raul M. Grijalva (AZ); Virginia Foxx (NC); Joe Courtney (CT); John A. Yarmuth (KY); Tim Walberg (MI); Glenn Thompson (PA); Frederica S. Wilson (FL); Suzanne Bonamici (OR); Mark Takano (CA); Joaquin Castro (TX); Mark Pocan (WI); Alma S. Adams (NC); Donald Norcross (NJ); Rick W. Allen (GA); Elise M. Stefanik (NY); Glenn Grothman (WI); James Comer (KY); Jim Banks (IN); Adriano Espaillat (NY); Pramila Jayapal (WA); Joseph D. Morelle (NY); Susan Wild (PA); Jahana Hayes (CT); Russ Fulcher (ID); Andy Levin (MI); Haley M. Stevens (MI); Ilhan Omar (MN); Mikie Sherrill (NJ); Fred Keller (PA); Gregory Murphy (NC); Michelle Steel (CA); Mariannette Miller-Meeks (IA); Mary E. Miller (IL); Frank J. Mrvan (IN); Victoria Spartz (IN); Kathy E. Manning (NC); Madison Cawthorn (NC); Teresa Leger Fernandez (NM); Jamaal Bowman (NY); Mondaire Jones (NY); Diana Harshbarger (TN); Burgess Owens (UT); Bob Good (VA); Scott Fitzgerald (WI); Kweisi Mfume (MD) |
| Witnesses | Marr, Linsey Ph.D., Professor of Civil and Environmental Engineering, Virginia Polytechnic Institute and State University, Blacksburg, VA; Michaels, David, Ph.D., Professor of Occupational and Environmental Medicine, The George Washington University, Former Assistant Secretary of OSHA, Washington, DC; Muhindura, Pascaline, RN, COVID Progressive Care Unit, Research Medical Center, on behalf of National Nurses United, Kansas City, MO; Rath, Manesh, Partner, Keller and Heckman LLP, Washington, DC |
| Serial Numbers | Serial No. 117-1 |
| Bill Numbers | H.R. 1180 |
| | H.R. 1195 |
| | H.R. 1319 |

---

```
[House Hearing, 117 Congress]
[From the U.S. Government Publishing Office]


            CLEARING THE AIR: SCIENCE-BASED STRATEGIES
                    TO PROTECT WORKERS FROM
                      COVID 19 INFECTIONS

=======================================================================

                             HEARING

                           BEFORE THE

                        SUBCOMMITTEE ON
                      WORKFORCE PROTECTIONS

                            OF THE

                 COMMITTEE ON EDUCATION AND LABOR
                 U.S. HOUSE OF REPRESENTATIVES

                 ONE HUNDRED SEVENTEENTH CONGRESS

                          FIRST SESSION

                          _____

          HEARING HELD IN WASHINGTON, DC, MARCH 11, 2021

                          _____

                      Serial No. 117-1

                          _____
```

# But:
# Hearings are not research ready

No search, nor filtered export
> Dates, topics
> Members of Congress, Witnesses
> Bills
> Individual speeches

Have to create the dataset ourselves!
1. Harvest it.
2. Parse it.
3. Store it.

# YOU WILL LEARN 3 TOOLS

1.  Harvesting the data:
    Headless **web scraping** using R.

2.  Parsing the data:
    Introduction to **Regular Expressions**.

3.  Organizing the data:
    Basics of **relational database design**.

All material at:
https://github.com/johannesfritz/scu-harvesting

# 1.
# HARVESTING THE DATA

Headless web scraping using R

# Two approaches to harvesting websites

1. GET requests
   > Effectively downloading the HTML/etc. files
   > R library {httr}


2. Browser automation + GET requests
   > Used for website testing
   > «headless»: without open window (GUI)
   > R library {webdriver}

# WHY WE NEED BROWSER AUTOMATION.

► Don't know all the URLs.

► Site is dynamically rendered.

# {WEBDRIVER} BASICS

►Library supporting «headless» browsing
>  Using phantomJS
>  Intuitive for rapid use
>  Flexible extensions (if you know JavaScript)
>  Alternative: {Rselenium}


►Supports
>  Browsing
>  Manipulating the page

# {WEBDRIVER} BASICS
## BROWSING

session = browser window

Log onto URL

like your browser icons

extract HTML or screenshot

```
s <- Session$new(host = "127.0.0.1", port = 8910)

s$delete()
s$status()


s$go(url)
s$getUrl()
s$goBack()
s$goForward()
s$refresh()
s$getTitle()
s$getSource()
s$takeScreenshot(file = NULL)
```

# {WEBDRIVER} BASICS
## LOCATING + MANIPULATING TARGETS

search the HTML source code →

extract text (not code) →

manipulate entry fields →

```
e$findElement(css = NULL, linkText = NULL,
        partialLinkText = NULL, xpath = NULL)
e$findElements(css = NULL, linkText = NULL,
        partialLinkText = NULL, xpath = NULL)

e$isSelected()
e$getValue()
e$setValue(value)
e$getAttribute(name)
e$getClass()
e$getCssValue(name)
e$getText()
e$getName()
e$getData(name)
e$getRect()
e$isEnabled()
e$click()
e$clear()
e$sendKeys(...)
```

# OUR INITIAL SCRIPT

```
1   library(webdriver)
2   library(XML)
3
4   ## starting the headless browser
5   my.browser = run_phantomjs()
6   my.browser
7
8
9   # Opening a session and accessing the URL
10  my.url="https://www.govinfo.gov/app/collection/chrg/"
11  my.session = Session$new(port= my.browser$port)
12
13  my.session$go(my.url)
14
15
16  # Looking around the page
17  ## verifying I arrived at the URL
18  my.session$getUrl()
19  my.session$getTitle()
20
21  ## Taking screenshots
22  ## Very useful to check you're navigating right
23  my.session$takeScreenshot()
24
```

# Super useful: Taking screenshots

# LOAD THE SCRIPT + SEE IF IT WORKS.
## TRY OUT YOUR OWN PAGE, TOO

[Copy it from GitHub](#)

# MOVING ON:
# KEY ELEMENTS OF OUR FULL SCRIPT

1. Navigate the site:
   need to tell our browser where to click
   > Tool: XPath

2. Create our local copy:
   store all pages of interest

3. Make capture comprehensive:
   go through all the pages we need
   > Tool: for loops

# 1. Navigate the site
## the Document Object model (DOM)

# HTML PRIMER
## TAGS + ATTRIBUTES

| Tag name | Function |
|---|---|
| `<head>` | The head of a page includes general information such as its authors, language etc. |
| `<body>` | All the content you see in the browser is inside the body tags. |
| `<p>` | A paragraph of text. |
| `<br>` | A break between two lines of text. |
| `<b>` , `<i>` , `<u>` | Display text in bold, italics or underlined. |
| `<a href="gpo.html">` | A link to the page gpo.html. |
| `<img src="gpo.jpg">` | Displays the picture gpo.jpg. |
| `<table>` , `<tr>` , `<td>` | A table including its rows (tr) and cells (td). |
| `<div>` | Divides an HTML document into different parts. Usually used for layout purposes. |
| `<span>` | Similar to `<div>` tag, but only for small parts of the code. |

attribute

# HTML PRIMER
## DOCUMENT STRUCTURE

S. Hrg. 113 - DEPARTMENT OF DEFENSE APPROPRIATIONS FOR FISCAL YEAR 2014

Appropriation. Wednesday, April 24, 2013.                    PDF | Text | More

```
<table class="browse-node-table">
<tr id="this-is-the-1st-row">
<td colspan="2" id="this-is-the-only-cell"> S. Hrg. 113 - DEPARTMENT OF DEFENSE
APPROPRIATIONS FOR FISCAL YEAR 2014 </td> </tr>
<tr id="this-is-the-2nd-row">
<td id="this-is-one-cell"> <span class="results-line2"> Appropriation. Wednesday, April 24,
2013. </span> </td>
 <td id="this-is-another-cell"> <a href="https://www.gpo.gov/link.to.PDF" target="_blank">
PDF </a> <a href="https://www.gpo.gov/link.to.text" target="_blank"> Text </a> <a
href="https://www.gpo.gov/link.to.more" target="_blank"> More </a> </td> </tr> <tr
id="this-is-the-3rd-row"> <td colspan="2" id="this-is-the-only-cell"> </td> </tr>
 </table>
```

# HTML PRIMER

## DOCUMENT STRUCTURE - INDENTED

```html
<table class="browse-node-table">

  <tr id="this-is-the-1st-row">

    <td colspan="2" id="this-is-the-only-cell">

    S. Hrg. 113 - DEPARTMENT OF DEFENSE APPROPRIATIONS FOR FISCAL YEAR 2014

    </td>

  </tr>

  <tr id="this-is-the-2nd-row">

    <td id="this-is-one-cell">

      <span class="results-line2">

        Appropriation. Wednesday, April 24, 2013.

      </span>

    </td>

    <td id="this-is-another-cell">

      <a href="https://www.gpo.gov/link.to.PDF" target="_blank">

        PDF

      </a>

      <a href="https://www.gpo.gov/link.to.text" target="_blank">

        Text
```
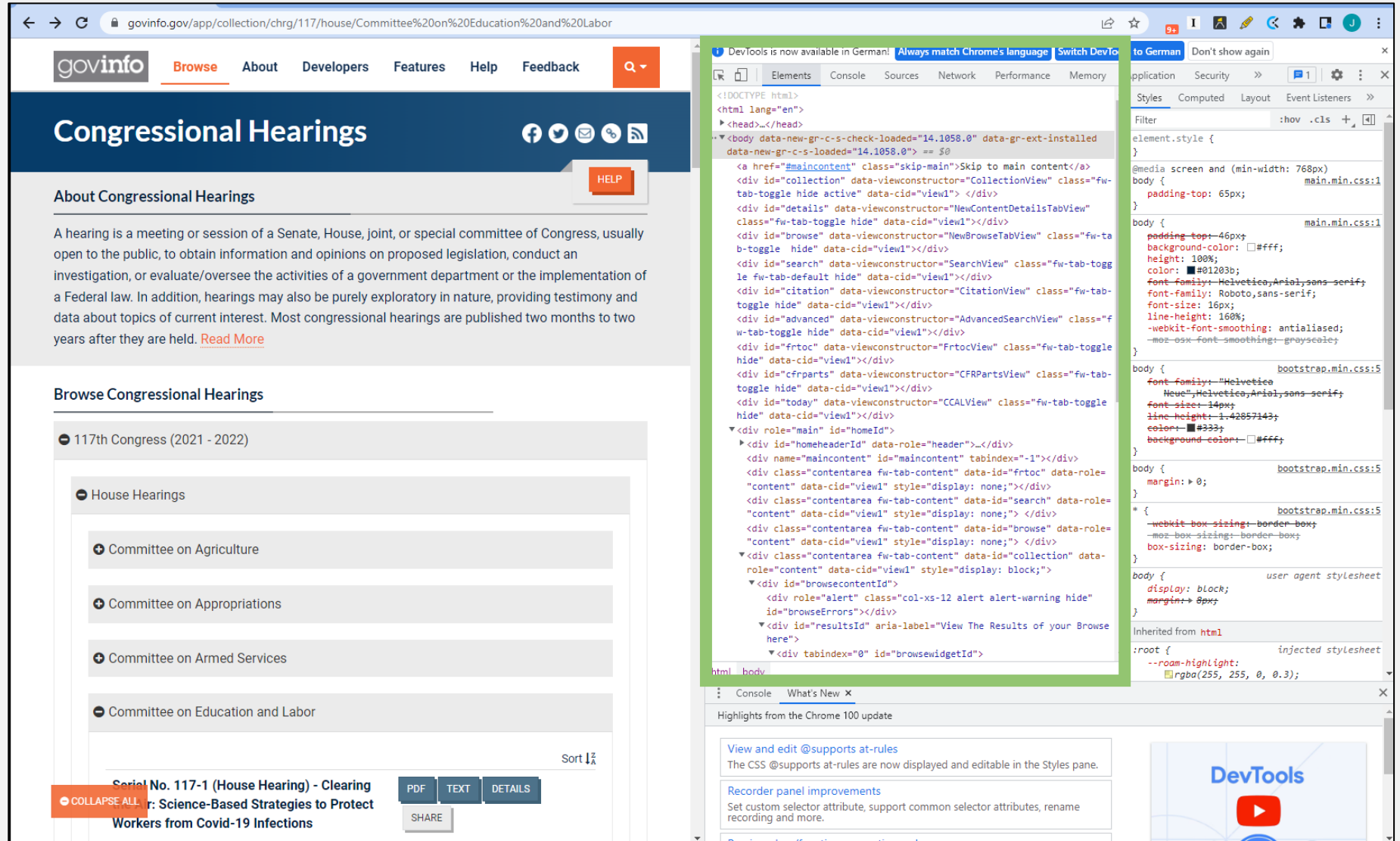
indentation levels

1

2

3

4

# XPath
## Navigating the DOM

► Simple query language to describe DOM locations.

► Basic rules:
1. We only include the opening tags. The program will just ignore closing tags it sees inside the code along the way.
2. The path only includes tags where the program has to move into a new nest (or towards the right in the above visualization).
3. We use numbers in squared brackets to identify a specific tag in cases where there are multiple ones of the same kind within the same nest.
4. We connect the different pieces of our path with backslashes in-between each tag.

# XPath example
## CLICKING ON THE «TEXT» LINK

```html
<table class="browse-node-table">

  <tr id="this-is-the-1st-row">

    <td colspan="2" id="this-is-the-only-cell">

    S. Hrg. 113 - DEPARTMENT OF DEFENSE APPROPRIATIONS FOR FISCAL YEAR 2014

    </td>

  </tr>

  <tr id="this-is-the-2nd-row">

    <td id="this-is-one-cell">

      <span class="results-line2">

        Appropriation. Wednesday, April 24, 2013.

      </span>

    </td>

    <td id="this-is-another-cell">

      <a href="https://www.gpo.gov/link.to.PDF" target="_blank">

        PDF

      </a>

      <a href="https://www.gpo.gov/link.to.text" target="_blank">

      Text
```

``table``          "Start at the beginning of our code (``table``),

        ``tr[2]``          use the second row tag (``tr[2]``).

          ``td[2]``          Within that row tag, go into the second cell tag you can find (``td[2]``), and

      ``a[2]``stop in that cell's second link tag (``a[2]``)."

➔ ``/table/tr[2]/td[2]/a[2]``

```html
<table class="browse-node-table">

    <tr id="this-is-the-1st-row">

        <td colspan="2" id="this-is-the-only-cell">

        S. Hrg. 113 - DEPARTMENT OF DEFENSE APPROPRIATIONS FOR FISCAL YEAR 2014

        </td>

    </tr>

    <tr id="this-is-the-2nd-row">

        <td id="this-is-one-cell">

          <span class="results-line2">

            Appropriation. Wednesday, April 24, 2013.

          </span>

        </td>

        <td id="this-is-another-cell">

          <a href="https://www.gpo.gov/link.to.PDF" target="_blank">

            PDF

          </a>

          <a href="https://www.gpo.gov/link.to.text" target="_blank">

          Text
```

# Finding the XPath
## using the «Inspector» in chrome/firefox

Open with:

► CTRL + SHIFT + I

► ⌥ + ⌘ + I

# Finding the XPath
## Select element

# FINDING THE XPATH
## COPY PREFERRED XPATH VERSION

# SHORT XPATHS PROTECT AGAINST HTML CODING ERRORS

► Full XPath
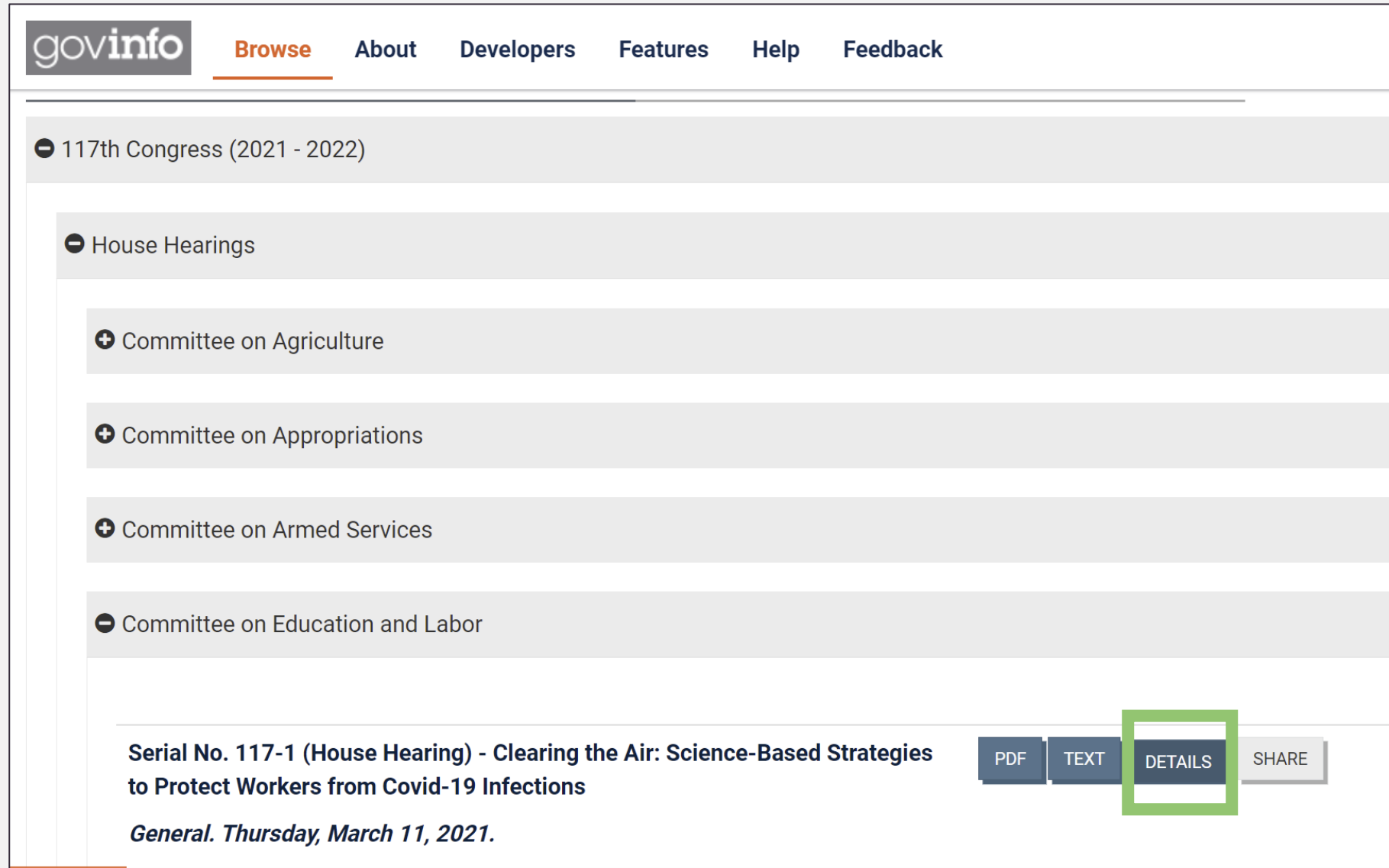/html/body/div[10]/div[6]/div/div[2]/div/div[1]/div[2]/div/div[2]/div/div[3]/p

► Short XPath
//*[@id="aboutDescription"]/p

► CSS selector (less flexible)
#aboutDescription > p

# Adding to our script: navigation + storage

```
53   ## extracting the full HTML file into the 'data' folder
54   html = my.session$getSource()
55
56   writeLines(text = html, con = 'data/my first page.html')
```

# 2. CREATE A LOCAL COPY
## CAN YOU CAPTURE THIS DETAILS PAGE?

govinfo

**Browse**   About   Developers   Features   Help   Feedback

⊖ 117th Congress (2021 - 2022)

⊖ House Hearings

⊕ Committee on Agriculture

⊕ Committee on Appropriations

⊕ Committee on Armed Services

⊖ Committee on Education and Labor

**Serial No. 117-1 (House Hearing) - Clearing the Air: Science-Based Strategies to Protect Workers from Covid-19 Infections**

*General. Thursday, March 11, 2021.*

PDF   TEXT   DETAILS   SHARE

check your result:
url

# Harvesting online data

Short course at Santa Clara University, ECON 173

Johannes Fritz

5/12/22

# YOU WILL LEARN 3 TOOLS

1.  Harvesting the data:
    Headless **web scraping** using R.

2.  Parsing the data:
    Introduction to **Regular Expressions**.

3.  Organizing the data:
    Basics of **relational database design**.

All material at:
https://github.com/johannesfritz/scu-harvesting

# Re-organised the code into the steps we took

# FROM LAST WEEK

```
1   library(webdriver)
2   library(XML)
3
4   ## starting the headless browser
5   my.browser = run_phantomjs()
6   my.browser
7
8
9   # Opening a session and accessing the URL
10  my.url="https://www.govinfo.gov/app/collection/chrg/"
11  my.session = Session$new(port= my.browser$port)
12
13  my.session$go(my.url)
14
15
16  ## Navigating on the page to the details page of a specific hearing and storing that page.
17
18  ## (1) Locating the first menu item '117th Congress (2021 - 2022)'
19  my.target= my.session$findElement(xpath='/html/body/div[10]/div[6]/div/div[3]/div/div/div[1]/div[1]/div')
20  my.target$getText() ## Seeing whether I got it
21  my.target$click()   ## Clicking on it to open the menu itself
```

# MORE XPATH: AXES

► If you have a prominent position in the code (e.g. an attribute).

► Axes allow you to navigate from there more flexibly without stating the path

► Axes names
  > parent or child
  > ancestor or descendant
  > preceding-sibling or following-sibling

# HTML PRIMER
## DOCUMENT STRUCTURE - INDENTED

**ancestor**

```
<table class="browse-node-table">

    <tr id="this-is-the-1st-row">

        <td colspan="2" id="this-is-the-only-cell">

        S. Hrg. 113 - DEPARTMENT OF DEFENSE APPROPRIATIONS FOR FISCAL YEAR 2014

        </td>

    </tr>
```

**parent**
**start**
**child**

```
    <tr id="this-is-the-2nd-row">

        <td id="this-is-one-cell">

            <span class="results-line2">

            Appropriation. Wednesday, April 24, 2013.

            </span>

        </td>
```

**following-sibling**

```
        <td id="this-is-another-cell">

            <a href="https://www.gpo.gov/link.to.PDF" target="_blank">

            PDF
```

//td[@id='this-is-one-cell']/following-sibling::td/a[1]

```
            </a>

            <a href="https://www.gpo.gov/link.to.text" target="_blank">

            Text
```

# XPath resource

**W3schools: XPath tutorial**
https://www.w3schools.com/xml/xpath_intro.asp

► Syntax

► Axes

► Special operators

# Moving on:
# Key elements of our full script

1. Navigate the site:
   need to tell our browser where to click
   > Tool: XPath

2. Create our local copy:
   store all pages of interest

3. Make capture comprehensive:
   go through all the pages we need
   > Tool: for loops

# 3. MAKE CAPTURE COMPREHENSIVE

# Code gets repetitive fast

```
16  ## Navigating on the page to the details page of a specific hearing and storing that page.
17
18  ## (1) Locating the first menu item '117th Congress (2021 - 2022)'
19  my.target= my.session$findElement(xpath='/html/body/div[10]/div[6]/div/div[3]/div/div/div[1]/div[1]/div')
20  my.target$getText() ## Seeing whether I got it
21  my.target$click()   ## Clicking on it to open the menu itself
22
23
24  ## (2) Locating the first menu item 'House Hearings'
25  my.target= my.session$findElement(xpath='/html/body/div[10]/div[6]/div/div[3]/div/div/div[1]/div[2]/div/div[1]/div[1]/div/span')
26  my.target$getText()
27  my.target$click()
28
29
30  ## (3) Locating the first menu item 'Committee on Education and Labor'
31  my.target= my.session$findElement(xpath='/html/body/div[10]/div[6]/div/div[3]/div/div/div[1]/div[2]/div/div[1]/div[2]/div/div[4]/div[1]/div/span')
32  my.target$getText()
33  my.target$click()
34
35
36  ## (4) Locating the first DETAILS button for the hearing on ' Clearing the Air: Science-Based Strategies to Protect Workers from Covid-19 Infections'
37  my.target= my.session$findElement(xpath='/html/body/div[10]/div[6]/div/div[3]/div/div/div[1]/div[2]/div/div[1]/div[2]/div/div[4]/div[2]/div/table[1]/tbody/tr/td[2]/div/a[3]')
38  my.target$getText()
39  my.target$click()
40
41  my.session$getUrl() ## did we get to the right location?
```

1. Move repetitions into a for loop
2. Extract loop elements directly from the HTML
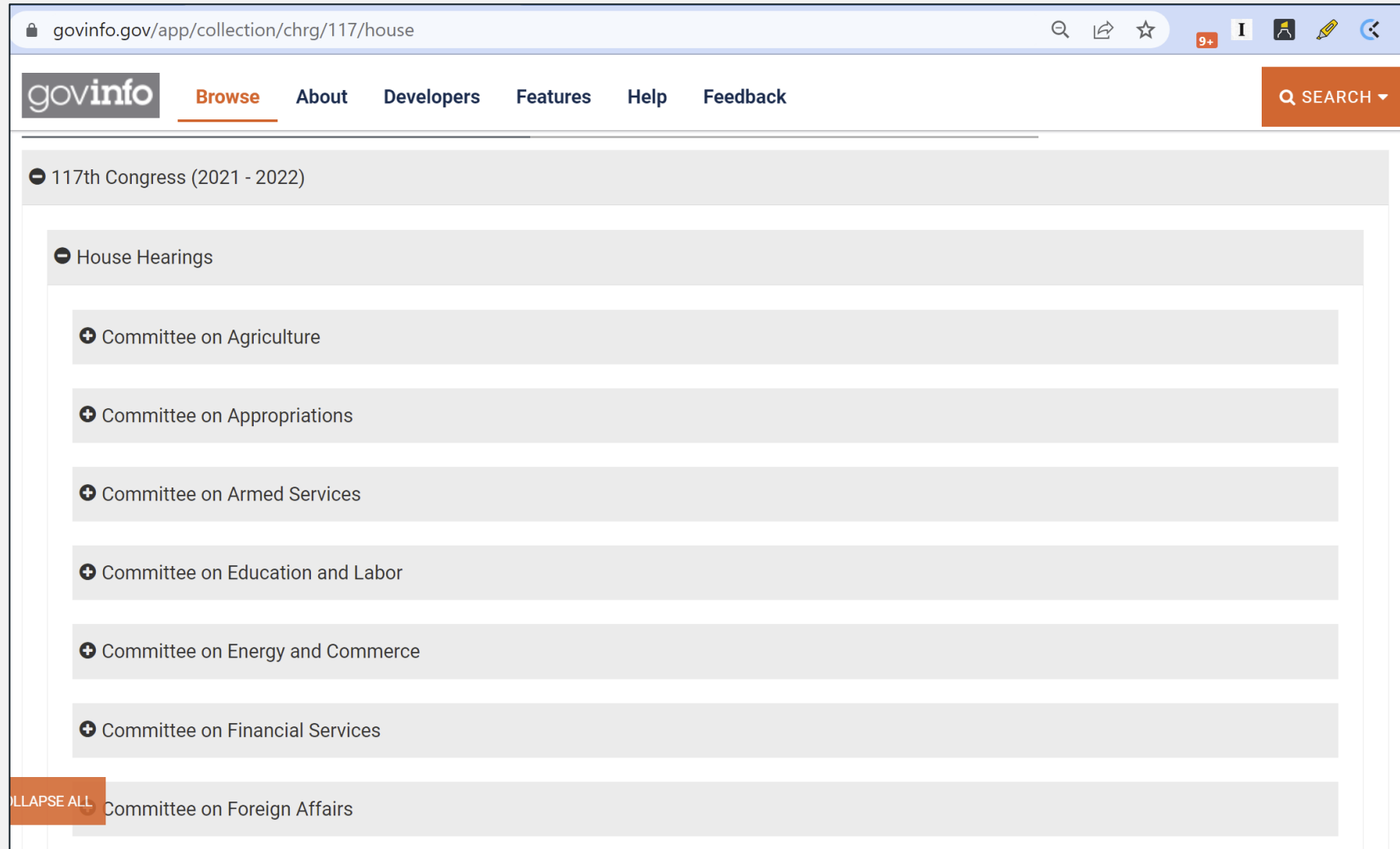
# FOR LOOPS
## COLLECTING FROM THOUSANDS OF PAGES

► You have to repeat the same operation x times

► for loops avoid code repetition

```
for (count in 1:1000) {
  print(count)
}
```

# MOVE REPETITIONS INTO A FOR LOOP

```
15
16   ## Navigating on the page: Same as in script 2 but now all in one go
17
18   i.want.to.click=c('/html/body/div[10]/div[6]/div/div[3]/div/div/div[1]/div[1]/div',
19                     '/html/body/div[10]/div[6]/div/div[3]/div/div/div[1]/div[2]/div/div[1]/div[1]/div/span',
20                     '/html/body/div[10]/div[6]/div/div[3]/div/div/div[1]/div[2]/div/div[1]/div[2]/div/div[4]/div[1]/div/span',
21                     '/html/body/div[10]/div[6]/div/div[3]/div/div/div[1]/div[2]/div/div[1]/div[2]/div/div[4]/div[2]/div/table[1]/tbody/tr/td[2]/div/a[3]')
22
23   for(this.element in i.want.to.click){
24     print(paste("Clicking on:",this.element))
25     my.target= my.session$findElement(xpath=this.element)
26     my.target$getText() ## Seeing whether I got it
27
28     ## Clicking on it to open the menu itself
29     my.target$click()
30     Sys.sleep(.5) ## need a pause of .5s because else the browser can't reload as fast as we are clicking
31
32   }
33
```

# Extract loop elements directly from the HTML

# CREATE A LIST OF ALL ELEMENTS IN AN XPATH LOCATION

```
<table class="browse-node-table">

  <tr id="this-is-the-1st-row">

    <td colspan="2" id="this-is-the-only-cell">

      S. Hrg. 113 - DEPARTMENT OF DEFENSE APPROPRIATIONS FOR FISCAL YEAR 2014

    </td>

  </tr>

  <tr id="this-is-the-2nd-row">

    <td id="this-is-one-cell">

      <span class="results-line2">

        Appropriation. Wednesday, April 24, 2013.

      </span>

    </td>

    <td id="this-is-another-cell">

      <a href="https://www.gpo.gov/link.to.PDF" target="_blank">

        PDF

      </a>

      <a href="https://www.gpo.gov/link.to.text" target="_blank">

        Text

      </a>

      <a href="https://www.gpo.gov/link.to.more" target="_blank">

        More

      </a>

    </td>

  </tr>
```
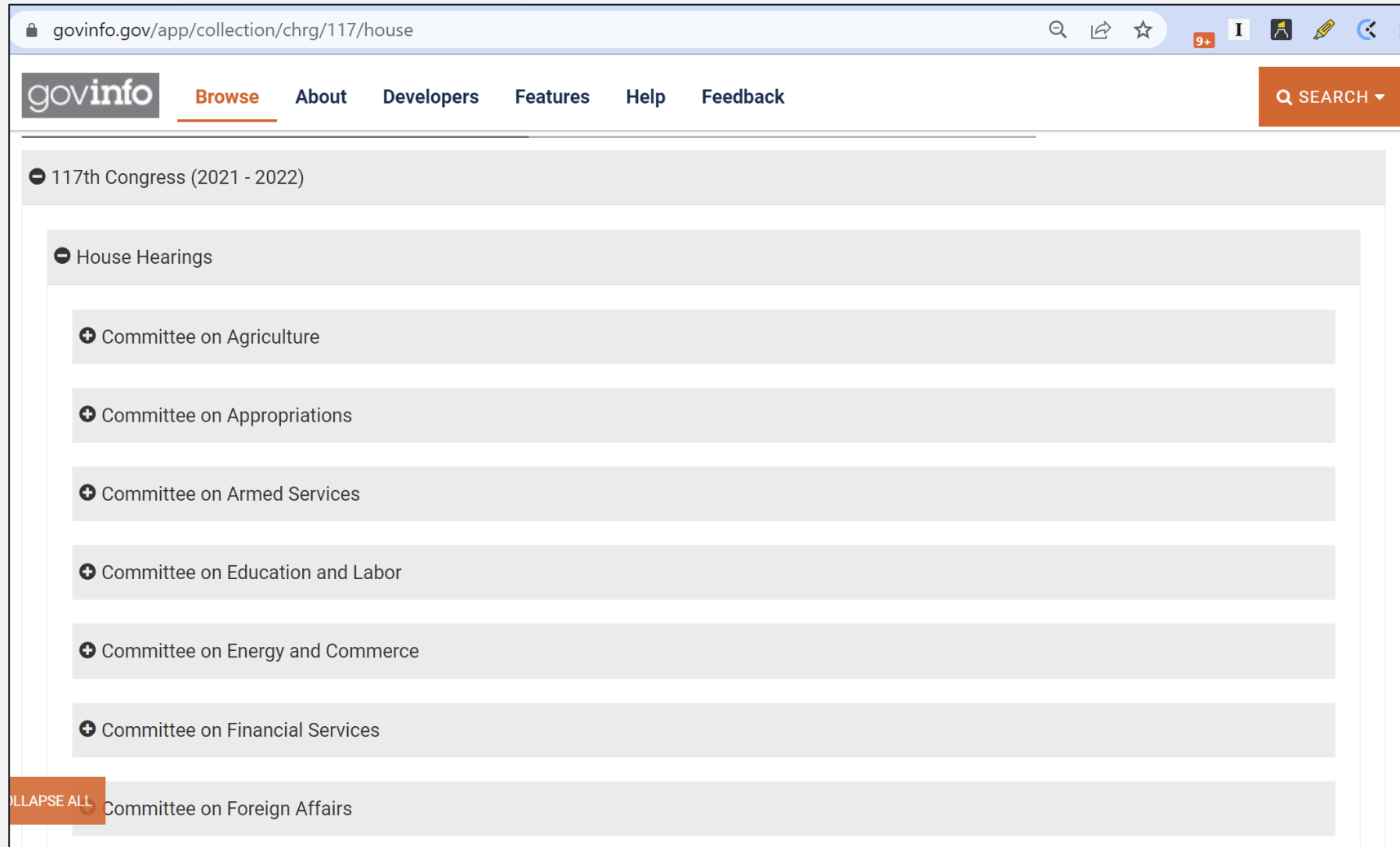
# EXTRACT THE COMMON ROOT

XPaths to all links in this example

> //table/tr[2]/td[2]/a[1]
> //table/tr[2]/td[2]/a[2]
> //table/tr[2]/td[2]/a[3]

common root

► Finding the common root yields 3 elements.

# FOR LOOP EXERCISE:
# PRINT THE NAMES OF ALL HOUSE COMMITTEES INTO YOUR R CONSOLE

govinfo.gov/app/collection/chrg/117/house

**govinfo**

Browse    About    Developers    Features    Help    Feedback

🔍 SEARCH ▾

⊖ 117th Congress (2021 - 2022)

⊖ House Hearings

⊕ Committee on Agriculture

⊕ Committee on Appropriations

⊕ Committee on Armed Services

⊕ Committee on Education and Labor

⊕ Committee on Energy and Commerce

⊕ Committee on Financial Services

COLLAPSE ALL    Committee on Foreign Affairs

# 2.
# PARSING THE DATA

Introduction to Regular Expressions.

# USUAL TEXT SEARCH IS VERY RIGID

# This is the query as received by the computer

**S** **Character.** Matches a "S" character (char code 83). Case sensitive.

**a** **Character.** Matches a "a" character (char code 97). Case sensitive.

**n** **Character.** Matches a "n" character (char code 110). Case sensitive.

**t** **Character.** Matches a "t" character (char code 116). Case sensitive.

**a** **Character.** Matches a "a" character (char code 97). Case sensitive.

**Character.** Matches a SPACE character (char code 32).

**C** **Character.** Matches a "C" character (char code 67). Case sensitive.

**l** **Character.** Matches a "l" character (char code 108). Case sensitive.

**a** **Character.** Matches a "a" character (char code 97). Case sensitive.

**r** **Character.** Matches a "r" character (char code 114). Case sensitive.

**a** **Character.** Matches a "a" character (char code 97). Case sensitive.

# RegEx flexibility helps you clean up + extract from text

# REGEX SUPPORTS FLEXIBLE SEARCH

► {Santa Clara est. 1851}

► Features to describe an expression:
> Character types
> Character quantifiers
> Character location + sequence

# ReGex commands

## CHARACTER TYPES

► **.** any kind of symbol

► **\d** all or **\D** no digits 0-9

► **\w** all or **\W** no alphanumeric characters (a-z, 0-9)

► **\s** all or **\S** no whitespace (as well as \n \t \v)

► **[rlCtnaS]** any of the characters in any order

> **[A-Z]** all capitalized letters

> **[a-z]** all letters in small type

> **[0-9]** all numbers

# ReGex commands
## Character quantifiers

► Default x=1

► * is x=0 or more

► + is x= 1 or more

► {n,m} is n<= x <=m

► ? means min(x) incl. zero x

► () designates groups

# ReGex commands
## Character locations, sequence + groups

► **^** means "start of string"
► **$** means "end of string"
► **(?=abc)** matches something written before "abc"
► **(?!abc)** matches something not written before "abc"
► **(?<=abc)** matches something written after "abc"
► **(?<!abc)** matches something not written after "abc"

# ReGex commands
## Avoiding special characters

► What do you do when you want to find a special RegEx character?
> `+ ? . $` etc.

► Add a backslash (in R: two):
> `\\+`
> `\\?`
> `\\.`
> `\\$`

# RegEx in R

► {base} library
> Matching:
  » grepl(pattern, x, ignore.case = FALSE)
  » grepl("", "") →

> Substitute:
  » gsub(pattern, replacement, x, ignore.case = FALSE)
  » gsub("", "", "") →

► {stringr} library
> Several additional functions.
  » Extraction:        str_extract_all(string, pattern)
  » Match count:   str_count(string, pattern)

# RegEx practice exercise

► Please open the practice script

# Best RegEx playground
## regexr.com

# REGEX RESOURCES

► RegEx sandbox website: https://regexr.com/

► Extensive documentation: http://www.regular-expressions.info

► Flexible text editor: https://www.sublimetext.com/

# 3.
# Organizing the data

Basics of relational database design

# Single-table setup quickly inefficient
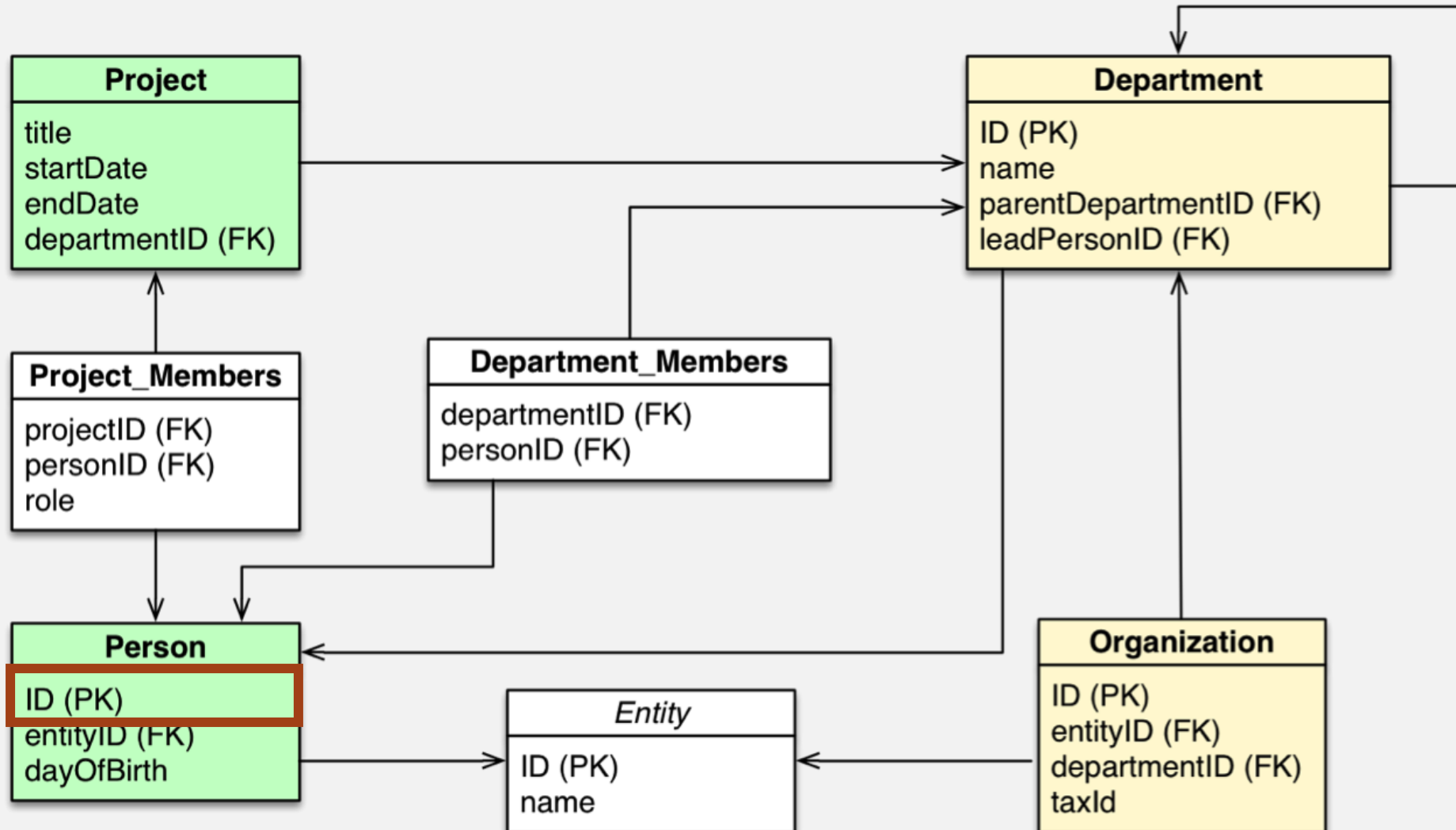
Sources

► Repetitions

► Unused variables

► Large volumes (of text)

Result:

► Lots of computational power

► Low speed

# BASIC SETUP FOR A LIST OF PROJECTS AND THEIR MEMBERS



**Project**
- title
- startDate
- endDate
- departmentID (FK)

**Department**
- ID (PK)
- name
- parentDepartmentID (FK)
- leadPersonID (FK)

**Project_Members**
- projectID (FK)
- personID (FK)
- role

**Department_Members**
- departmentID (FK)
- personID (FK)

**Person**
- ID (PK)
- entityID (FK)
- dayOfBirth

**Entity**
- ID (PK)
- name

**Organization**
- ID (PK)
- entityID (FK)
- departmentID (FK)
- taxId

ID's to organise and merge the data

# THAT WAS IT …

1. Harvesting the data:
   Headless **web scraping** using R.

2. Parsing the data:
   Introduction to **Regular Expressions**.

3. Organizing the data:
   Basics of **relational database design**.

All material at:
https://github.com/johannesfritz/scu-harvesting

# Thank you!

END OF SHORT COURSE