

Topic 10: Web Scraping, String Manipulation, and Regular Expressions



Outline

- Web scraping
- String variables and regular expressions
- Scraping the Williams website

Web Pages Are Built in **H**yper**T**ext **M**arkup **L**anguage (HTML)

ECON 370



Pamela lakiela sollabors

This is the website for Professor Pamela Inkiela's ECON 370 course at Williams College Data Science for Economic Analysis

Data Science for Economics

Course Description:

relevant for economic analysis including data visualization, marbine use these data science tools differently than many researchers in statistics prediction. Through a combination of lectures, hands-on labs, and group analyze economic data using modern data science techniques in R or Python

Course Information:

Syllabus Schedule

Projects:

Data Visualization Project

```
2 chtml langu"en-US">
      costs charest-"UTE-0">
      cmeta http-equiv="X-UA-Compatible" content="IF-edge">
      cmeta_name="viewnort"_content="width-device-width, initial-scale=1">
 | </-- Begin Jekyll SEO tag v2.8.0 -->
 ctttle>Data Science for Economics | ECON 370c/tttle>
10 <meta name-"generator" content-"Jekyll v3.10.0" />
11 cmeta property-"og:title" content-"Data Science for Economics" />
12 cmeta property: "or:locale" content: en US" />
1) costs name "description" costent-"course materials for data science for economic analysis" /a
14 <meta property="og:description" content="course materials for data science for economic analysis" />
15 clink rel="canonical" href="https://pjakiela.github.io/ECON370/" />
15 ceets property-"origin" contest-"https://piakiela.github.in/6008378/" />
17 (mate property, "origina piec" contents "ECON 170" ()
is (meta property-"og:type" content-"website" />
10 ceets names"twitterscand" contents"summary" />
10 cents property."twitten:title" content."Oats Science for Economics" /s
31 country types"annlication(ldsison")
22 ("@context": "https://schema.org", "@type": "website", "description": "course materials for data science for economic
  analysis", "headline", "Data Science for Economics", "name", "ECON 378", "nublisher"; ("Styne"; "Organization", "loro";
   ("Stype": "ImageObject", "url": "https://piakiela.github.io/ECOME79/economist-data-crop-vz.jpg"), "url": "https://
  piakiela.github.io/ECON370/*)c/script>
23 cl -- End Jehull SEO ton -->
      clink rel="stylesheet" href="/FCON178/assets/css/style.css/v=a3f172edcf2af85e5418b475cb5c99928484cafd">
      vaccint arc-"https://cdeis.cloudflore.com/clou/fibs/htmlSchiu/3.7.3/htmlSchiu.min.de"sv/scripts
      <!!fendifl-->
    e Chands
    chades
      edia elass-"womones">
          chiacdiv style="text-align: right"aca href="https://plakiela.github.in/propurs/"accom 370c/aac/divac/hta
            cime sec="/FCON370/economist.data.cron.v2.ing" alt="logo" />
          cdiv style="text-align: right">ccmall>source: The Economist c/small>c/div>
          coacdiv style="text-align: right">Instructor: c/diva
          cdiv style="text-align: right"sca bref="https://piakiela.github.io/"spamela Jakiela c/asc/divs
          cp>cdiv style="text-align: right">ca href="https://piakiela.github.io/ECON178/">home </div>
          cdiv style="text-align: right">ca href="https://pjakiela.github.jo/ECON370/ECON370-
  syllabus-2024-09-11.pdf">syllabus </a></div>
```

Elements and Attributes

A web page is made up of **elements** that are arranged in a hierarchical structure

- Elements start and end with a tag: for example, <title>ECON 370</title>
 - Every html page contains an html element (<html>...</html>)
 - ► The html element contains the elements (children) head and body
- Many elements appear multiple times within a page: for example, a paragraph
- Tags can contain attributes that encode element-specific information, for example:
 - ► <h1 id="data-science-for-economics">Data Science for Economics</h1>
 - Syllabus
- Sometimes data is stored as a table within the page: ...

When to Scrape

Simple web scraping tools are useful when:

- Data is stored in a table
- A page contains many repetitions of the same structure/sequence of elements, and you want to combine those elements into a data frame for analysis

```
Example: <h4><span class="course_code">ECON 105</span><span
class="course_terms">(F)</span> <span class="course_code">SEM</span> <span
class="course_title">Gender in the Global Economy</span></h4><h4><span
class="course_code">ECON 107</span><span class="course_terms_blank"></span>
<span class="course_code">SEM</span> <span class="course_title">Inequality in a
Classless Society: The Soviet Experiment and its Aftermath</span></h4>
```

- A sequence of (ideally numerically indexed) pages use the same structure and elements, and you want to combine the information from multiple pages to build a data set
 - **Example:** building a data set of all recent NBER working papers

How to Scrape: Get the HTML, Extract Text and Attributes

- 1. Do the actual scraping: load HTML into R or Python
- 2. Extract elements and the text and attributes they contain:
 - Extract elements with tag "a": html_elements(mypage, "a") or mypage.select("a")
 - Extract elements with class equal to a with (".a")
 - Extract elements with id equal to a with ("#a")
 - Extraction tasks can be sequenced (first extract parents, then specific children)
 - SelectorGadget can (sometimes) make this process easier
 - Final step is to extract text (printed on web page) or the value of an attribute

How to Scrape: Example

https://ephsports.williams.edu/sports/womens-soccer/roster

When Not to Scrape

Many websites cannot be scraped easily with simple tools (can you do it?)

Policy advice: don't try to scrape websites that don't want to be scraped

Web scraping raises a range of ethical issues (should you do it?)

- Some personally-identifiable information posted on the web should not be used for research
 - In general, work products and things shared under creative comments licenses are in-bounds, personal (non-professional) content posted with some expectation of privacy is out-of bounds
 - ▶ When collecting personally identifiable information for research, check with your IRB
 - Collecting personally-identifiable information is often unnecessary
- Also: don't violate a websites terms of service (if you are bound by them) or copyright law
- Finally, be polite: to avoid over-burdening a server, always build in pauses between queries

String Variables

A string variable is a variable that is stored as a sequence of characters

- Categorical variables can be stored as strings or as labeled/indexed numeric variables
 - Relatively small range of possible valid values
 - Multiple observations with the same value
- Text data that does not represent a (small-ish) set of mutually exclusive categories
 - Examples: names, phone numbers, open-ended responses

The distinction between text variables and categorical variables can be subtle, context-specific

- "What is your name?" vs. "Who is your favorite Beatle?"
- "Pamela Jakiela" vs. "Ann Miller"

String Manipulation

- Split them into multiple strings
 - ► Split name into firstname and lastname using the delimiter ","
 - ► "Harrison, George" → "Harrison" + "George"
- Combine multiple strings
 - Combine lastname and firstname into fullname using the delimiter ""
 - "Harrison" + "George" → "George Harrison"
- Detect or count (substring) matches
 - Does a phone number contain the area code 413?
 - How many items did a customer purchase?
- Replace substrings

String Variables

How might we answer: does a phone number contain the area code 413?

- Split phone_number into area code, the rest of the phone number
 - ▶ A good approach when generating a categorical variable for metropolitan area of residence, but maybe not if you just want to know if someone is from western Massachusetts
- Check whether phone_number contains the substring "413"
 - Consider the phone number: 510-408-1413

We actually want to know if phone_number starts with 413

Regular Expressions

A **regular expression** or **regex** is a string combining (normal) characters and **metacharacters** that can be parsed and then used to identify and manipulate pattern matches in text data

• Example: "^413" means starts with 413

Regular expression tools are useful for:

- Expanding the scope of pattern matching
 - "[Pp]rofessor [Jj]akiela", "[0-9]", "[0-9]*", "[0-9]+"
- Restricting the scope of pattern matching
 - "^Professor", "Jakiela\$", "[a-z]{2,}ing"

To match a character that is also a metacharacter (e.g. "\$"), use a slash: "\\$"

• To match a "\", you need to use "\\"

Lab #10

Lab #10 can be completed in groups of up to 4 people, one submission per group

• Your group needs to write a script that scrapes data from Williams websites and processes it in order to answer (up to) 4 questions about sports teams and academic departments

The ECON370-web-williams-soccer-example scripts illustrate how to scrape a data table or other specific elements from simple web pages in either R (rvest) or Python (BeautifulSoup)

- Adapt the script to scrape and process data using these libraries
- Feel free to use SelectorGadget to identify the appropriate selectors