

Cologne Center for Comparative Politics

A Primer to Web Scraping with R

Introduction

Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

Welcome

Welcome!

Thank you for your interest in the course!

Welcome!

Thank you for your interest in the course!

I'm glad that I will have the opportunity to give you a gentle introduction to web scraping with R.

Welcome!

Thank you for your interest in the course!

I'm glad that I will have the opportunity to give you a gentle introduction to web scraping with R.

Just a few words on my personal background:

- Lecturer in Political Data Science at the Hertie School of Governance
- political scientist by training
- working with web-based data since about 2010
- what I do with web data: measure public awareness, news consumption, political behavior

Your background and interests

“Please let me know if there is any topic that you definitely want me to cover.”

Your background and interests

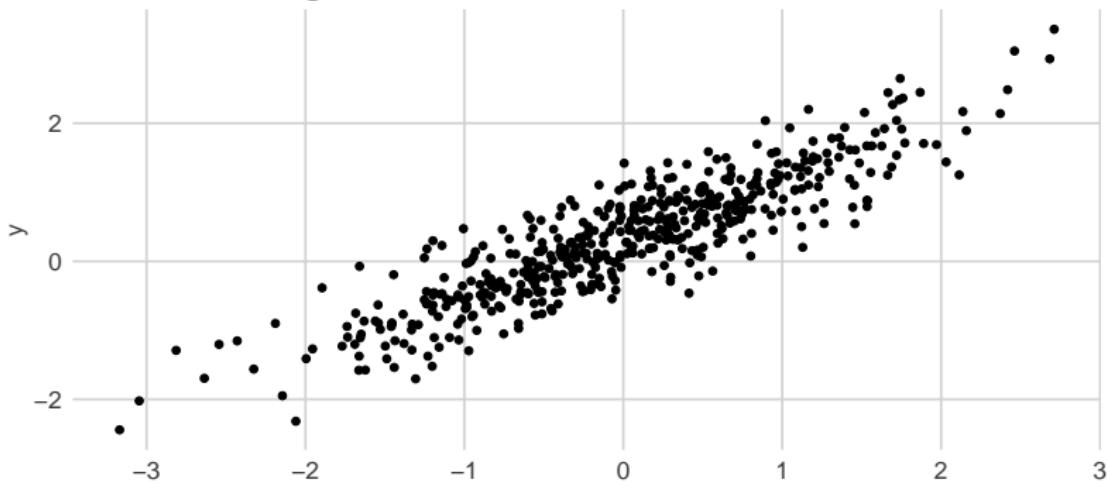
“Please let me know if there is any topic that you definitely want me to cover.”

- "Scraping Flash pages."
- "rather interested in general introduction"
- "How to scrape Java scripts"
- "I am just getting acquainted with R studio"
- "linear regression"

I want all of you to be happy, so...

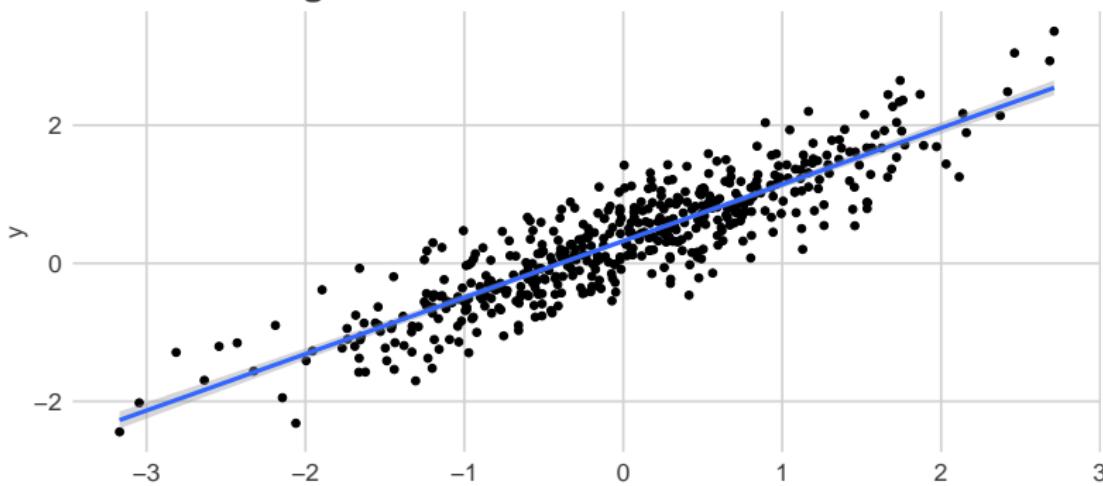
I want all of you to be happy, so...

How linear regression works.



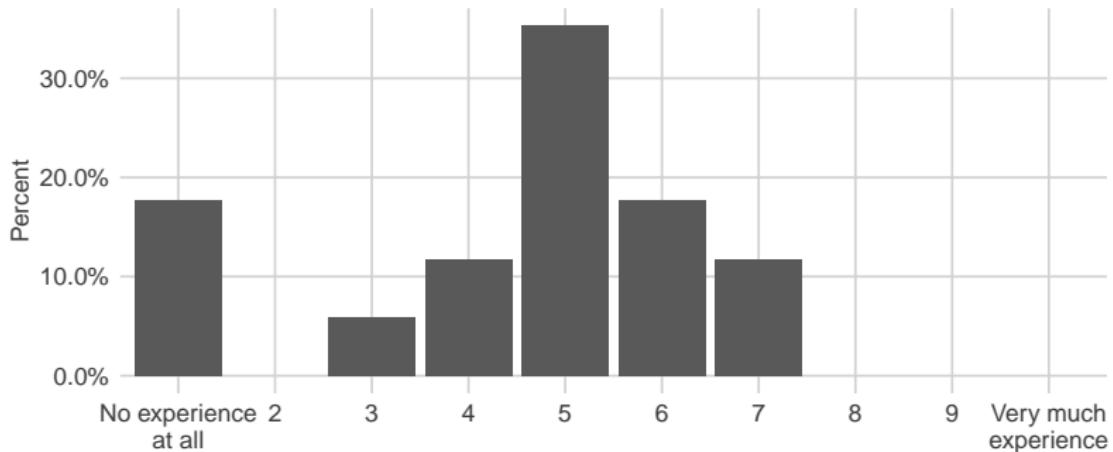
I want all of you to be happy, so...

How linear regression works.



Your background and interests

"Please rate your experience with R!"



Source: Participant Poll

Your background and interests



Source: Participant Poll

Goals and outline

Goals

After attending this course, ...

Goals

After attending this course, ...

- you have acquired basic knowledge of web technologies

Goals

After attending this course, ...

- you have acquired basic knowledge of web technologies
- you are able to scrape information from static and dynamic websites using R

Goals

After attending this course, ...

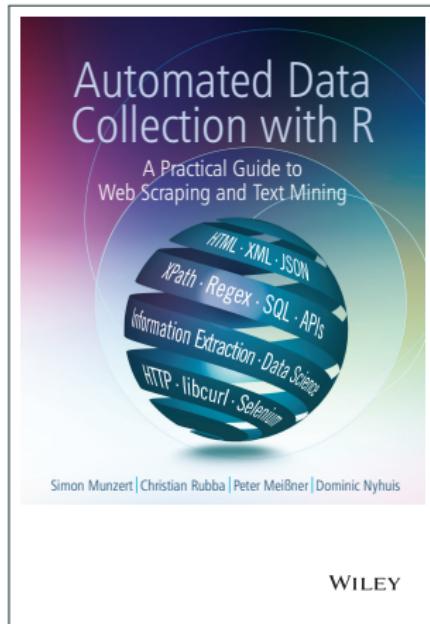
- you have acquired basic knowledge of web technologies
- you are able to scrape information from static and dynamic websites using R
- you are able to access web services (APIs) with R

Course outline

Time	Topic
09:00 - 09:45	Introduction; a first encounter with the Web using R
09:45 - 12:15	Scraping static webpages
12:15 - 13:15	<i>Lunch break</i>
13:15 - 14:00	Scraping dynamic webpages
14:00 - 15:30	Tapping APIs
15:30 - 16:30	Scraping ethics and workflow

The accompanying book

- contains most of which I tell you during the course (but much more, and at times more accurate)
- written between 2012 and 2014 → not entirely up-to-date anymore, but I will provide updated material during the course
- homepage with materials: www.r-datacollection.com
- if you find any errors in the book, please let me know!



WILEY

Cologne Center for Comparative Politics

A Primer to Web Scraping with R

Overview

Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

Web scraping with R

Web scraping. What? Why?

Web scraping

A.k.a. screen scraping, is the business of

- getting (unstructured) data from the web and
- bringing it into shape (e.g., clean, make tabular format)

Web scraping. What? Why?

Web scraping

A.k.a. screen scraping, is the business of

- getting (unstructured) data from the web and
- bringing it into shape (e.g., clean, make tabular format)

A data analyst's view

- data abundance online
- social interaction online
- services track social behavior

Web scraping. What? Why?

Web scraping

A.k.a. screen scraping, is the business of

- getting (unstructured) data from the web and
- bringing it into shape (e.g., clean, make tabular format)

A data analyst's view

- data abundance online
- social interaction online
- services track social behavior

A pragmatist's view

- financial resources
- time resources
- reproducibility
- updateability

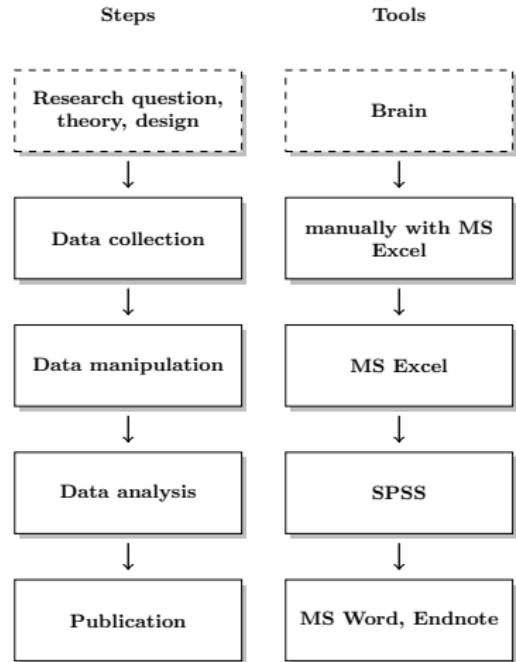
Why R?

Why R?

- free
- open source
- large community
- powerful tools for statistical analysis
- powerful tools for visualization
- flexible in processing all kinds of data/languages
- useful in every step of the workflow

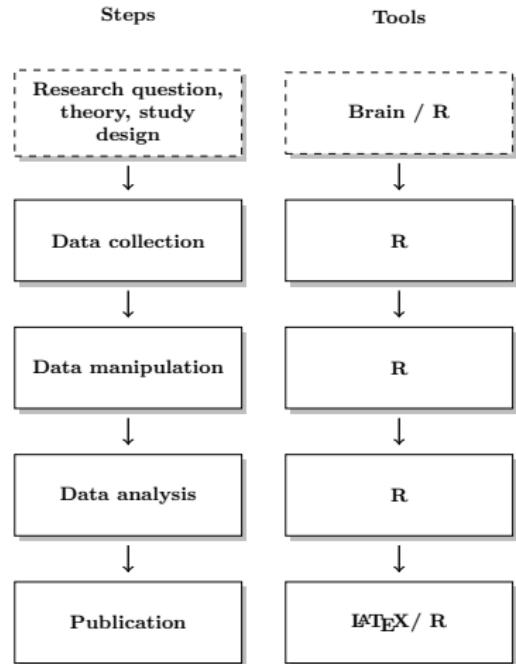
Why R?

- free
- open source
- large community
- powerful tools for statistical analysis
- powerful tools for visualization
- flexible in processing all kinds of data/languages
- useful in every step of the workflow



Why R?

- free
- open source
- large community
- powerful tools for statistical analysis
- powerful tools for visualization
- flexible in processing all kinds of data/languages
- useful in every step of the workflow

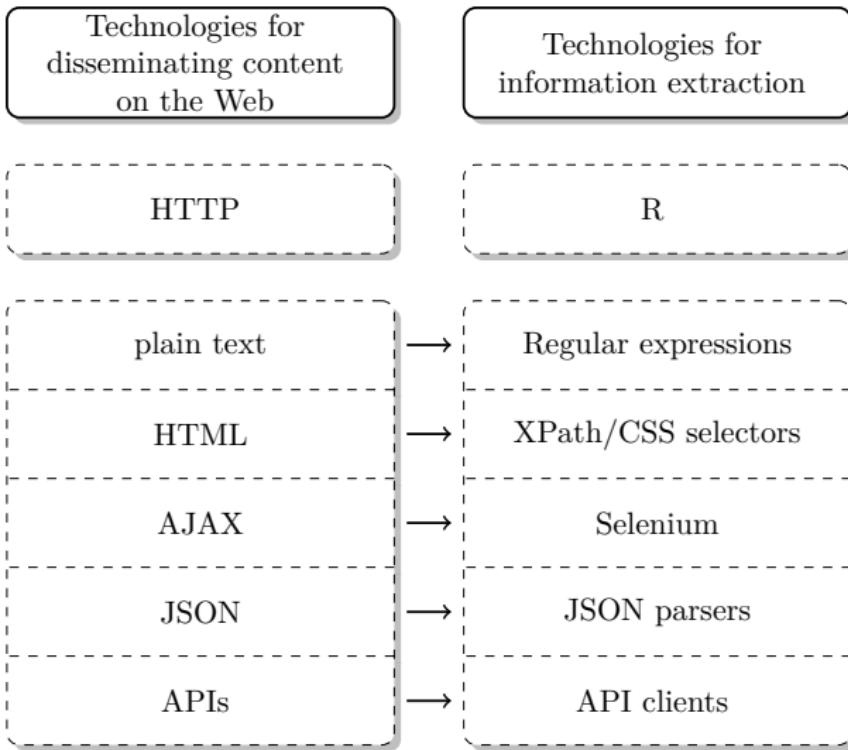


The philosophy behind web data collection with R

- no point-and-click procedure
- script the entire process from start to finish
- automation of
 - downloading
 - classical screen scraping
 - tapping APIs
 - parsing
 - data tidying, text data processing
- scaling up scraping procedures
- scheduling of scraping tasks

Technologies of the World Wide Web

Technologies of the World Wide Web



Technical setup

1. make sure that the newest version of R is installed on your computer (available here: <https://cran.r-project.org/>)
2. install the newest stable version of *RStudio Desktop* (available here: <https://www.rstudio.com/products/rstudio/download>)

Cologne Center for Comparative Politics

A Primer to Web Scraping with R

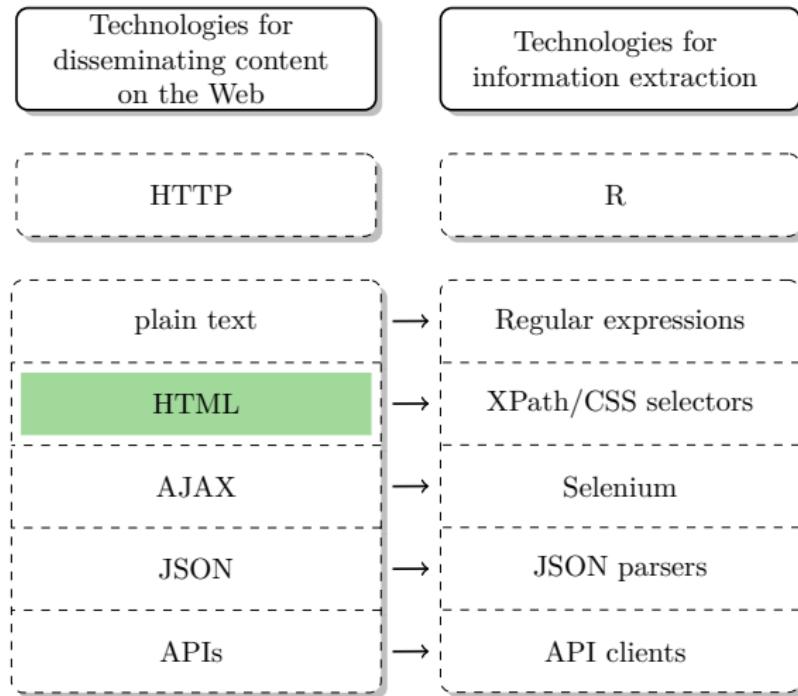
HTML

Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

Technologies of the World Wide Web



HTML basics

What's HTML?

- **HyperText Markup Language**
- markup language = plain text + markups
- W3C standard for the construction of websites
- lies underneath of what you see in your browser

HTML – a quick primer

What's HTML?

- HyperText Markup Language
- markup language = plain text + markups
- W3C standard for the construction of websites
- lies underneath of what you see in your browser

Why is this important to us?

- it determines where and how information is stored
- a basic understanding of HTML helps us locate the information we want to retrieve
- relax. A passive understanding of HTML is sufficient

HTML in the wild

Sicher <https://en.wikipedia.org/wiki/Berlin>

Not logged in Talk Contributions Create account Log in

Berlin

From Wikipedia, the free encyclopedia

This article is about the capital of Germany. For other uses, see Berlin (disambiguation).

Berlin (/bərˈlɪn/, German: [bɛʁˈliːn] (listen)) is the capital and the largest city of Germany as well as one of its constituent 16 states. With a population of approximately 3.5 million people,^[4] Berlin is the second most populous city proper and the seventh most populous urban area in the European Union.^[5] Located in northeastern Germany on the banks of rivers Spree and Havel, it is the centre of the Berlin-Brandenburg Metropolitan Region, which has about 6 million residents from more than 180 nations.^{[6][7][8][9]} Due to its location in the European Plain, Berlin is influenced by a temperate seasonal climate. Around one-third of the city's area is composed of forests, parks, gardens, rivers and lakes.^[10]

First documented in the 13th century and situated at the crossing of two important historic trade routes,^[11] Berlin became the capital of the Margravate of Brandenburg (1417–1701), the Kingdom of Prussia (1701–1918), the German Empire (1871–1918), the Weimar Republic (1919–1933) and the Third Reich (1933–1945).^[12] Berlin in the 1920s was the third largest municipality in the world.^[13] After World War II and its consequent occupation by the victorious countries, the city was divided; East Berlin became the capital of East Germany while West Berlin became a de facto West German exclave, surrounded by the Berlin Wall (1961–1989).^[14] Following German reunification in 1990, Berlin once again became the capital of Germany.



HTML in the wild

Berlin's mayors make up the council of mayors (*rat der Bürgermeister*), which is led by the city's Governing Mayor and advises the Senate. The neighborhoods have no local government bodies.

Twin towns – sister cities [edit]

See also: List of twin towns and sister cities in Germany

Berlin maintains official partnerships with 17 cities.^[100] Town twinning between Berlin and other cities began with its sister city Los Angeles in 1967. East Berlin's partnerships were canceled at the time of German reunification but later partially reestablished. West Berlin's partnerships had previously been restricted to the borough level. During the Cold War era, the partnerships had reflected the different power blocs, with West Berlin partnering with capitals in the Western World, and East Berlin mostly partnering with cities from the Warsaw Pact and its allies.

There are several joint projects with many other cities, such as Beirut, Belgrade, São Paulo, Copenhagen, Helsinki, Johannesburg, Mumbai, Oslo, Shanghai, Seoul, Sofia, Sydney, New York City and Vienna. Berlin participates in international city associations such as the Union of the Capitals of the European Union, Eurocities, Network of European Cities of Culture, Metropolis, Summit Conference of the World's Major Cities, and Conference of the World's Capital Cities. Berlin's official sister cities are:^[100]

- 1967 Los Angeles, United States
- 1987 Paris, France
- 1988 Madrid, Spain
- 1989 Istanbul, Turkey
- 1991 Warsaw, Poland^[101]
- 1991 Moscow, Russia
- 1992 Brussels, Belgium
- 1992 Budapest, Hungary^[102]
- 1993 Tashkent, Uzbekistan
- 1993 Mexico City, Mexico
- 1993 Jakarta, Indonesia
- 1994 Beijing, China
- 1994 Tokyo, Japan
- 1994 Buenos Aires, Argentina
- 1995 Prague, Czech Republic^[103]
- 2000 Windhoek, Namibia
- 2000 London, United Kingdom

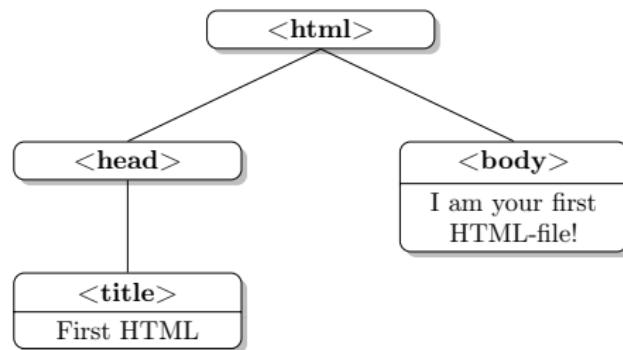
Capital city [edit]

Berlin is the capital of the Federal Republic of Germany. The President of Germany, whose functions are mainly ceremonial under the German constitution, has his official residence in Schloss Bellevue.^[104] Berlin is the seat of the German executive, housed in the Chancellery, the *Bundeskanzleramt*. Facing the Chancellery is the *Bundestag*, the German Parliament, housed in the renovated Reichstag building since the government relocated to Berlin in 1998. The *Bundesrat* ("federal council", performing the function of an upper house) is the representation of the Federal States (*Bundesländer*) of Germany and has its

HTML in the wild

Tree structure

```
1 <!DOCTYPE html>
2 <html>
3   <head>
4     <title id=1>First HTML</title>
5   </head>
6   <body>
7     I am your first HTML file!
8   </body>
9 </html>
```



Elements and attributes

Elements and attributes

Elements

Elements are a combination of *start tags*, content, and *end tags*.

Example:

```
1 <title>First HTML</title>
```

Syntax

element title	title
start tag	<title>
end tag	</title>
value	First HTML

Elements and attributes

Attributes

Attributes describe elements and are stored in the start tag. In HTML, there are specific attributes for specific elements.

Example:

```
1 <a href="http://www.r-datacollection.com/">Link to Homepage</a>
```

Syntax

- name-value pairs: `name="value"`
- simple and double quotation marks possible
- several attributes per element possible

Why tags and attributes are important

- tags structure HTML documents
- everything that structures a document can be used to extract information
- in the following, we get to know some important tags which are useful when scraping information from the Web

Important tags and attributes

Anchor tag <a>

- links to other pages or resources
- classical links are always formatted with an anchor tag
- the **href** attribute determines the target location
- the value is the name of the link

Link to another resource:

```
1 <a href="en.wikipedia.org/wiki/List_of_lists_of_lists">Link with absolute path</a>
```

Reference in a document:

```
1 <a id="top">Reference Point</a>
```

Link to a reference:

```
1 <a href="#top">Link to Reference Point</a>
```

Important tags and attributes

Heading tags `<h1>`, `<h2>`, ..., and paragraph tag `<p>`

- structure text and paragraphs
- heading tags range from level 1 to 6
- paragraph tag induces line break

Examples:

```
1 <p>This text is going to be a paragraph one day and separated from other  
2 text by line breaks.</p>
```

```
1 <h1>heading of level 1 -- this will be BIG</h1>  
2 ...  
3 <h6>heading of level 6 -- the smallest heading</h6>
```

Important tags and attributes

Listing tags , and <dl>

- the `` tag creates a numeric list, `` an unnumbered list, `<dl>` a definition list
- list elements are indicated with the `` tag

Example:

```
1 <ul>
2   ^^I<li>Dogs</li>
3   ^^I<li>Cats</li>
4   ^^I<li>Fish</li>
5 </ul>
```

Important tags and attributes

Organizational tags `<div>` and ``

- grouping of content over lines (`<div>`) or within lines (``)
- do not change the layout themselves but work together with CSS

Example of CSS definition

```
1  div.happy { color:pink;  
2      font-family:"Comic Sans MS";  
3      font-size:120% }  
4  span.happy { color:pink;  
5      font-family:"Comic Sans MS";  
6      font-size:120% }
```

In the HTML document

```
1  <div class="happy"><p>I am a happy styled paragraph</p></div>  
2  non-happy text with <span class="happy">some happiness</span>
```

Important tags and attributes

Form tag <form>

- allows to incorporate HTML forms
- client can send information to the HTTP server via forms
- whenever you type something into a field or click on radio buttons in your browser, you are interacting with forms

Example:

```
1 <form name="submitPW" action="Passed.html" method="get">
2   password:
3   <input name="pw" type="text" value="">
4   <input type="submit" value="SubmitButtonText">
5 </form>
```

Important tags and attributes

Table tags `<table>`, `<tr>`, `<td>`, and `<th>`

- standard HTML tables always follow a standard architecture
- the different tags allow to define the table as a whole, individual rows (including the heading), and cells
- if the data is hidden in tables, scraping will be straightforward

Example:

```
1 <table>
2   <tr> <th>Rank</th> <th>Nominal GDP</th> <th>Name</th> </tr>
3   <tr> <th></th> <th>(per capita, USD)</th> <th></th> </tr>
4   <tr> <td>1</td> <td>170,373</td> <td>Lichtenstein</td> </tr>
5   <tr> <td>2</td> <td>167,021</td> <td>Monaco</td> </tr>
6   <tr> <td>3</td> <td>115,377</td> <td>Luxembourg</td> </tr>
7   <tr> <td>4</td> <td>98,565</td> <td>Norway</td> </tr>
8   <tr> <td>5</td> <td>92,682</td> <td>Qatar</td> </tr>
9 </table>
```

Summary

Summary

- HTML is the *lingua franca* on the web
- content on webpages is structured by HTML tags that are nested in a tree structure
- to break open information, we will have to locate it in the HTML tree
- for web scraping purposes, a mostly passive knowledge of HTML is sufficient

```
<DIV>Q: HOW DO YOU ANNOY A WEB DEVELOPER?</SPAN>
```

Source: <https://xkcd.com/1144/> (Randall Munroe)

Cologne Center for Comparative Politics

A Primer to Web Scraping with R

XPath, Part I

Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

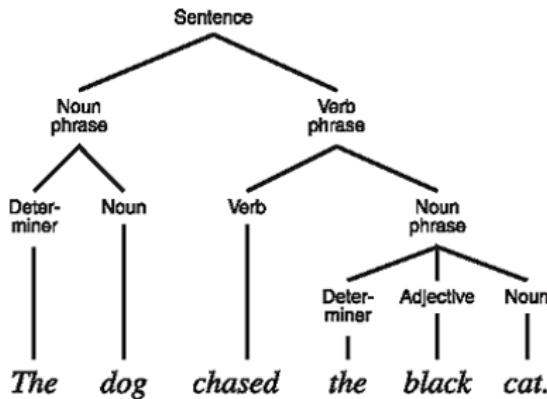
Accessing the HTML tree with R

Accessing the HTML tree with R

- HTML documents are human-readable
- HTML tags structure the document
- **web user perspective:** the browser interprets the code and renders the page
- **web scraper perspective:** use the tags to locate information; document has to be parsed first

Parsing

Parsing originally describes the syntactic analysis of text according to grammatical rules; analysis of the relationship between single parts of text. In programming, the input has to be interpreted (e.g., by R) to process the command.



Tools

- the `xml2` package allows us to parse XML-style documents
- the `rvest` package, which we will mainly use for scraping, wraps the `xml2` package, so we rarely have to load it manually
- HTML is a "flavor" of XML, so we can use the package to parse HTML
- one high-level function: `read_html()`
- `read_html()` represents the HTML in a list-style fashion
- we could also import HTML files via `readLines()`, but this is not parsing—the document's structure is not retained

HTML parsing with R

Parsing a website is straightforward

R code

```
1 library(rvest)  
2 parsed_doc <- read_html("https://google.com")
```

end

Functions to inspect the parsed document - better use the browser instead

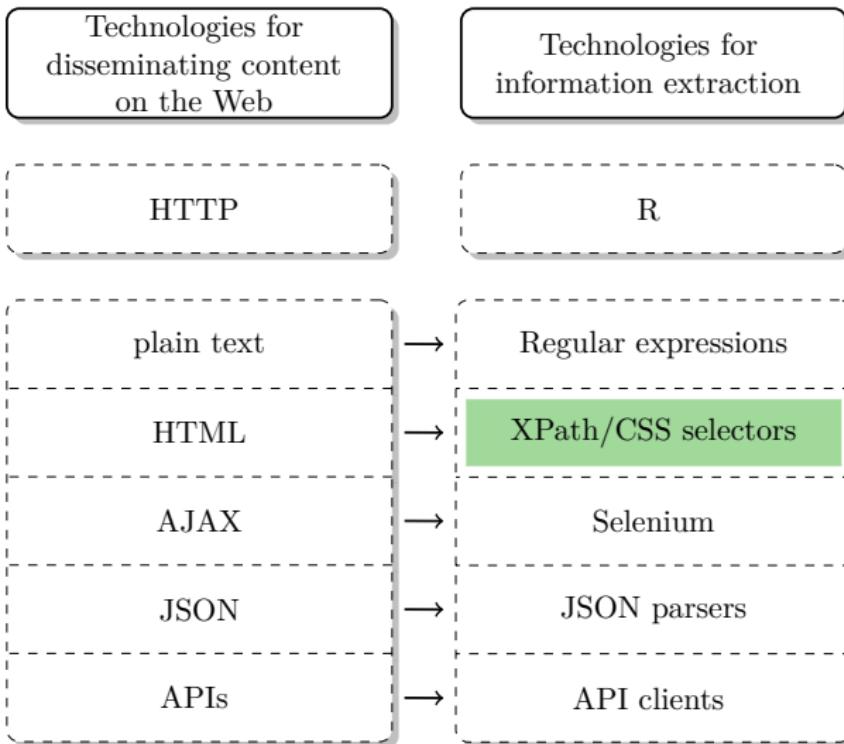
R code

```
3 html_structure(parsed_doc)  
4 as_list(parsed_doc)
```

end

XPath

Technologies of the World Wide Web



What's XPath?

Definition

- XML Path language, a W3C standard
- query language for XML-based documents (→ HTML)
- access node sets and extract content

What's XPath?

Definition

- XML Path language, a W3C standard
- query language for XML-based documents (→ HTML)
- access node sets and extract content

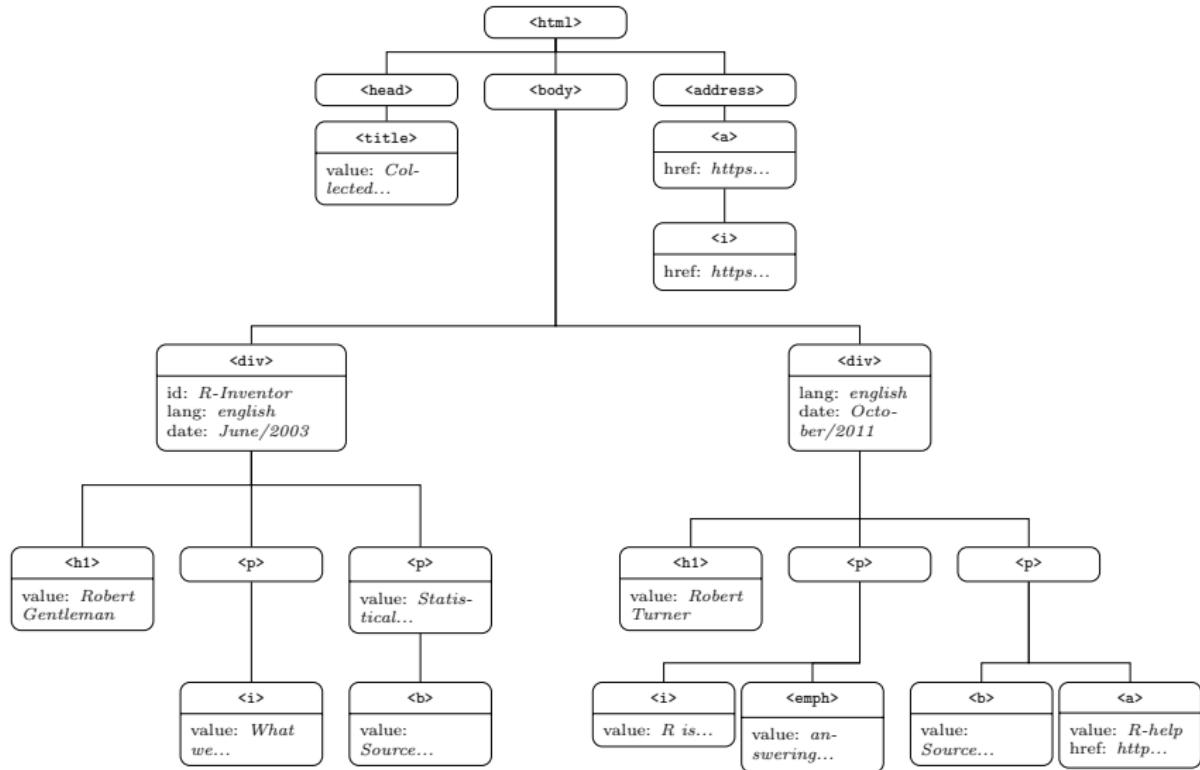
Why XPath for web scraping?

- source code of webpages (HTML) structures both layout and content
- not only content, but context matters!
- enables us to extract content based on its location in the document and (usually) regardless of its shape

Example

```
1 <!DOCTYPE HTML PUBLIC "-//IETF//DTD HTML//EN">
2 <html> <head>
3 <title>Collected R wisdoms</title>
4 </head>
5 <body>
6 <div id="R Inventor" lang="english" date="June/2003">
7   <h1>Robert Gentleman</h1>
8   <p><i>'What we have is nice, but we need something very different'</i></p>
9   <p><b>Source: </b>Statistical Computing 2003, Reisenburg</p>
10 </div>
11 <div lang="english" date="October/2011">
12   <h1>Rolf Turner</h1>
13   <p><i>'R is wonderful, but it cannot work magic'</i> <br><emph>answering a
14     request for automatic generation of 'data from a known mean and 95% CI'</
15     emph></p>
16   <p><b>Source: </b><a href="https://stat.ethz.ch/mailman/listinfo/r-help">R-help</
17     a></p>
18 </div>
19 </body>
20 <address><a href="http://www.r-datacollection.com"><i>The book homepage</i><a/></
21   address>
22 </html>
```

Example



Applying an XPath expression in R

- load package `rvest`
- parse document with `read_html()`
- query document with XPath expression using `html_nodes()`
- `rvest` can process XPath queries as well as CSS selectors
- in this course, we'll focus on XPath

R code

```
5 library(rvest)
6 parsed_doc <- read_html("../materials/fortunes.html")
7 html_nodes(parsed_doc, xpath = "//div[last()]/p/i")
{xml_nodeset (1)}
[1] <i>'R is wonderful, but it cannot work magic'</i>
```

end

The grammar of XPath

Grammar of XPath

Basic rules

1. we access nodes by writing down the hierarchical structure in the DOM that locates the node set of interest
2. a sequence of nodes is separated by slash symbols
3. the easiest localization of a node is given by the absolute path (but often not the most efficient one!)
4. apply XPath on document in R with the `html_nodes()` function

R code

```
8 html_nodes(parsed_doc, xpath = "//div[last()]/p/i")
{xml_nodeset (1)}
[1] <i>'R is wonderful, but it cannot work magic'</i>
```

end

Absolute vs. relative paths

- absolute paths start at the root node and follow the whole way down to the target node (with simple slashes, '/')
- relative paths skip nodes (with double slashes, '//')

R code

```
9 html_nodes(parsed_doc, xpath = "/html/body/div/p/i")
{xml_nodeset (2)}
[1] <i>'What we have is nice, but we need something very different'</i>
[2] <i>'R is wonderful, but it cannot work magic'</i>

10 html_nodes(parsed_doc, xpath = "//body//p/i")
{xml_nodeset (2)}
[1] <i>'What we have is nice, but we need something very different'</i>
[2] <i>'R is wonderful, but it cannot work magic'</i>
```

end

When to use absolute, when relative paths?

- relative paths faster to write
- relative paths often more comprehensive (but less robust)
- relative paths consume more computing time, as the whole tree has to be parsed, but this is usually of less relevance for reasonably small documents

R code

```
11 html_nodes(parsed_doc, xpath = "//i")
{xml_nodeset (3)}
[1] <i>'What we have is nice, but we need something very different'</i>
[2] <i>'R is wonderful, but it cannot work magic'</i>
[3] <i>The book homepage</i>
```

end

Wildcard operator

- meta symbol *
- matches any node
- works only for one arbitrary node
- far less important than wildcards in regular expressions

R code

```
12 html_nodes(parsed_doc, xpath = "/html/body/div/*/i")
  {xml_nodeset (2)}
[1] <i>'What we have is nice, but we need something very different'</i>
[2] <i>'R is wonderful, but it cannot work magic'</i>
13 # this does not work:
14 html_nodes(parsed_doc, xpath = "/html/body/*/i")
  {xml_nodeset (0)}
```

end

Navigational operators ‘.’ and ‘..’

- . accesses nodes at the same level ('self axis')
- useful when working with predicates
- .. accesses nodes at a higher hierarchical level

R code

```
15 html_nodes(parsed_doc, xpath = "//title/..")
{xml_nodeset (1)}
[1] <head>\n<meta http-equiv="Content-Type" content="text/html; charset=
...

```

end

Pipe operator

- combines several paths

R code —

```
16 html_nodes(parsed_doc, xpath = "//address | //title")
{xml_nodeset (2)}
[1] <title>Collected R wisdoms</title>
[2] <address>\n<a href="http://www.r-datacollection.com"><i>The book hom
...
————— end
```

Summary

Summary

- XPath is a little language that lets you query specific parts of an XML-style document
- it has its own grammar (logic) and vocabulary
- in this session, you learned the basics of XPath
- in the next session, you will learn more advanced, powerful XPath expressions



Source: https://commons.wikimedia.org/wiki/File:Mozie_Law_path_junction_-_geograph.org.uk_-_1131.jpg (Andy Stephenson)

Cologne Center for Comparative Politics

A Primer to Web Scraping with R

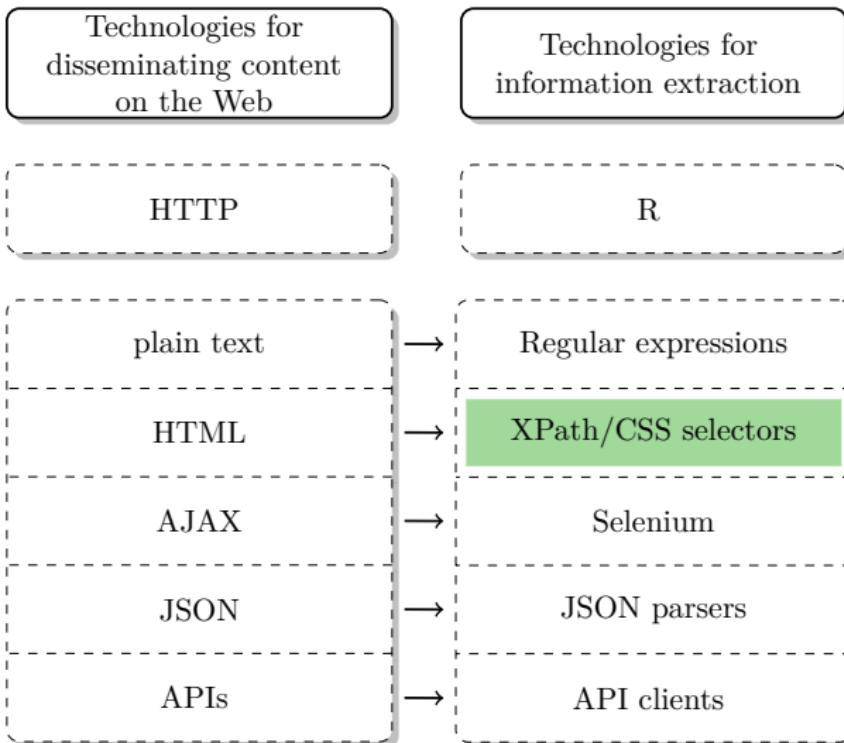
XPath, Part II

Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

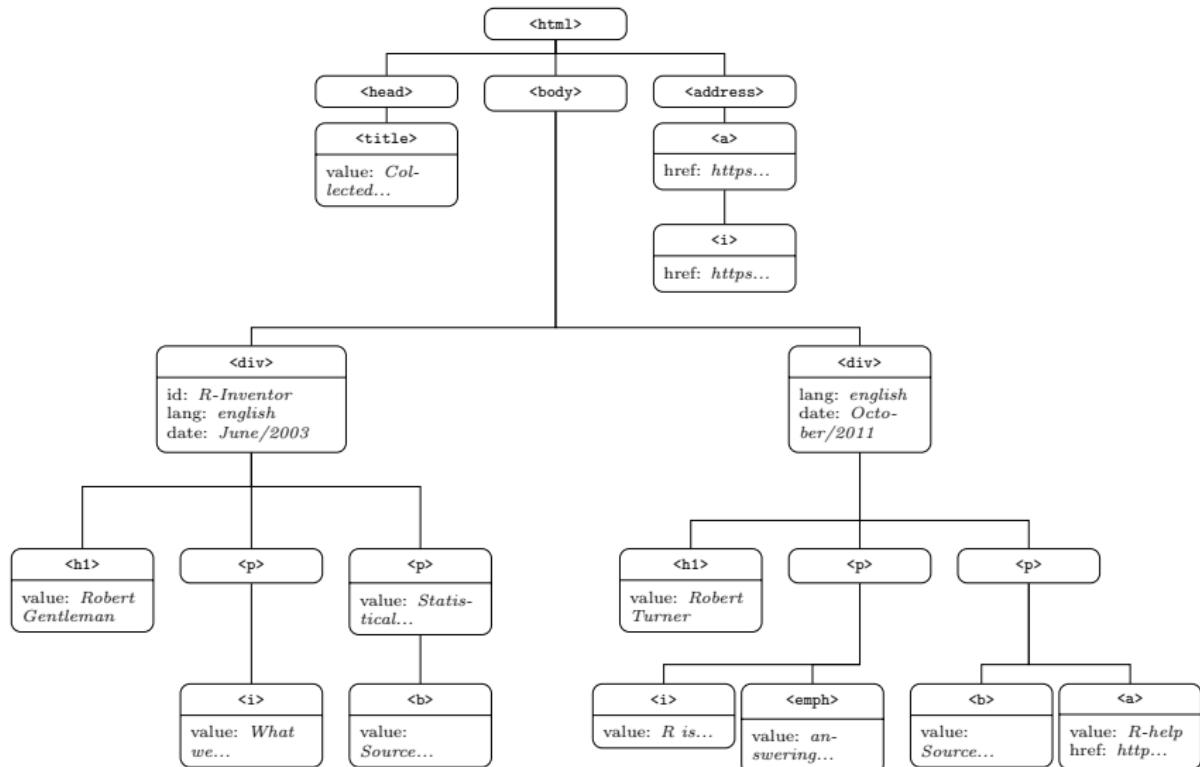
Technologies of the World Wide Web



Example

```
1 <!DOCTYPE HTML PUBLIC "-//IETF//DTD HTML//EN">
2 <html> <head>
3 <title>Collected R wisdoms</title>
4 </head>
5 <body>
6 <div id="R Inventor" lang="english" date="June/2003">
7   <h1>Robert Gentleman</h1>
8   <p><i>'What we have is nice, but we need something very different'</i></p>
9   <p><b>Source: </b>Statistical Computing 2003, Reisenburg</p>
10 </div>
11 <div lang="english" date="October/2011">
12   <h1>Rolf Turner</h1>
13   <p><i>'R is wonderful, but it cannot work magic'</i> <br><emph>answering a
14     request for automatic generation of 'data from a known mean and 95% CI'</
15     emph></p>
16   <p><b>Source: </b><a href="https://stat.ethz.ch/mailman/listinfo/r-help">R-help</
17     a></p>
18 </div>
19 </body>
20 <address><a href="http://www.r-datacollection.com"><i>The book homepage</i><a/></
21   address>
22 </html>
```

Example

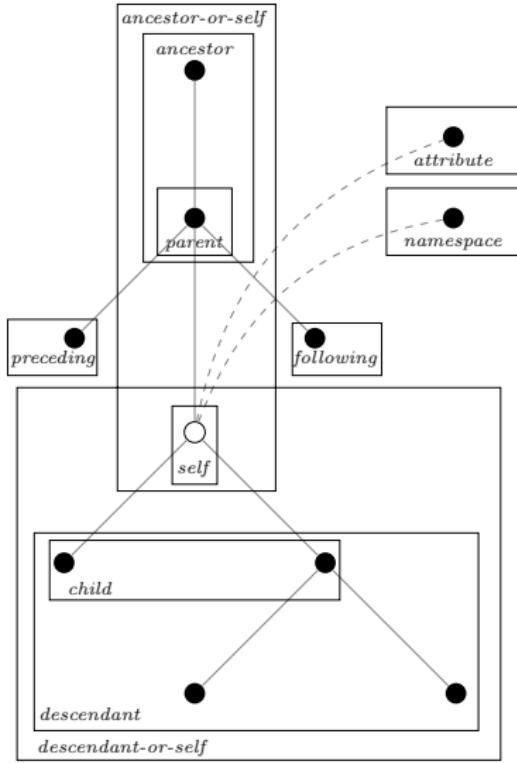


Node relations in XPath

'Family relations' between nodes

- the tools learned so far are sometimes not sufficient to access specific nodes without accessing other, undesired nodes as well
- relationship statuses are useful to establish unambiguity
- can be combined with other elements of the grammar
- basic syntax: `node1/relation::node2`
- we describe *relation* of `node2` to `node1`
- `node2` is to be extracted—we **always** extract the node at the end

Node relations in XPath



Node relations in XPath

Axis name	Result
ancestor	all ancestors (parent, grandparent, etc.) of the current node
ancestor-or-self	all ancestors of the current node and the current node itself
attribute	all attributes of the current node
child	all children of the current node
descendant	all descendants (children, grandchildren, etc.) of the current node
descendant-or-self	all descendants of the current node and the current node itself
following	everything in the document after the closing tag of the current node
following-sibling	all siblings after the current node
namespace	all namespace nodes of the current node
parent	the parent of the current node
preceding	all nodes that appear before the current node in the document, except ancestors, attribute nodes and namespace nodes
preceding-sibling	all siblings before the current node
self	the current node

Node relations in XPath

Example: access the `<div>` nodes that are ancestors to an `<a>` node:

R code

```
1 html_nodes(parsed_doc, xpath = "//a/ancestor::div")
{xml_nodeset (1)}
[1] <div lang="english" date="October/2011">\n  <h1>Rolf Turner</h1>\n  ...

```

end

Node relations in XPath

Example: access the `<div>` nodes that are ancestors to an `<a>` node:

R code

```
3 html_nodes(parsed_doc, xpath = "//a/ancestor::div")
{xml_nodeset (1)}
[1] <div lang="english" date="October/2011">\n    <h1>Rolf Turner</h1>\n    ...

```

end

Another example: Select all `<h1>` nodes that precede a `<p>` node:

R code

```
4 html_nodes(parsed_doc, xpath = "//p/preceding-sibling::h1")
{xml_nodeset (2)}
[1] <h1>Robert Gentleman</h1>
[2] <h1>Rolf Turner</h1>

```

end

Predicates

Predicates

- Conditions based on a node's features (`true/false`)
- applicable to a variety of features: name, value, attribute
- basic syntax:

node[predicate]

Commands in predicates

Function	Returns...
<code>name(<node>)</code>	name of <code><node></code> or the first node in a node set
<code>text(<node>)</code>	value of <code><node></code> or the first node in a node set
<code>@attribute</code>	value of a node's <i>attribute</i>
<code>string-length(str1)</code>	length of <code>str1</code> . If there is no string argument, it length of the string value of the current node
<code>translate(str1, str2, str3)</code>	<code>str1</code> by replacing the characters in <code>str2</code> with the characters in <code>str3</code>
<code>contains(str1,str2)</code>	TRUE if <code>str1</code> contains <code>str2</code> , otherwise FALSE
<code>starts-with(str1,str2)</code>	TRUE if <code>str1</code> starts with <code>str2</code> , otherwise FALSE
<code>substring-before(str1,str2)</code>	start of <code>str1</code> before <code>str2</code> occurs in it
<code>substring-after(str1,str2)</code>	remainder of <code>str1</code> after <code>str2</code> occurs in it
<code>not(arg)</code>	TRUE if the Boolean value is FALSE, and FALSE if the boolean value is TRUE
<code>local-name(<node>)</code>	name of the current <code><node></code> or the first node in a node set – without the namespace prefix
<code>count(<node>)</code>	count of a nodeset <code><node></code>
<code>position(<node>)</code>	index position of <code><node></code> that is processed
<code>last()</code>	number of items in the processed node list <code><node></code>

Predicates

Numeric predicates indicate positions, counts, etc.

Example: Select all first `<p>` nodes that are children of a `<div>` node:

R code

```
5 html_nodes(parsed_doc, xpath = "//div/p[position()=1]")
{xml_nodeSet (2)}
[1] <p><i>'What we have is nice, but we need something very different'</i>
...
[2] <p><i>'R is wonderful, but it cannot work magic'</i> <br><emph>answe
...

```

end

Further examples:

R code

```
6 html_nodes(parsed_doc, xpath = "//div/p[last()-1]")
7 html_nodes(parsed_doc, xpath = "//div[count(./*)>2]")
8 html_nodes(parsed_doc, xpath = "//*[string-length(text())>50]")

```

end

Textual predicates

- describe text features
- applicable on: node name, content, attributes, attribute values
- XPath 2.0 supports (simplified) regex, however, R (the **rvest** package and the underlying **xml2** package) does not support it

Example: Select all `<div>` nodes that contain an attribute named '`October/2011`':

R code

```
9 html_nodes(parsed_doc, xpath = "//div[@date='October/2011'])  
{xml_nodeset (1)}  
[1] <div lang="english" date="October/2011">\n    <h1>Rolf Turner</h1>\n    ...
```

end

Partial matching

- rudimentary string matching
- 'content contains' (`contains()`), 'content begins with' (`starts-with()`), 'content ends with' (`ends-with()`), 'content contains after split' (`substring-after()`)

R code

```
10 html_nodes(parsed_doc, xpath = "//*[contains(text(), 'magic')]")
{xml_nodeset (1)}
[1] <i>'R is wonderful, but it cannot work magic'</i>
```

end

Further examples:

R code

```
11 html_nodes(parsed_doc, xpath = "//div[starts-with(./@id, 'R')]")
12 html_nodes(parsed_doc, xpath = "//div[substring-after(./@date, '/') 
='2003']//i")
```

end

Content extraction

Extraction of text and attributes

Extraction

- until now: query of complete nodes
- common scenario: only parts of the node are interesting, e.g., content (value)
- additional extraction operations from a selected node set possible with additional extractor functions

Function	Argument	Return value
<code>html_text</code>		node value
<code>html_attr</code>	<code>name</code>	node attribute
<code>html_attrs</code>		(all) node attributes
<code>html_name</code>	<code>trim</code>	node name
<code>html_children</code>		node children

Extraction of node elements

Values/text

R code

```
13 html_nodes(parsed_doc, xpath = "//title") %>% html_text()  
[1] "Collected R wisdoms"
```

end

Attributes

R code

```
14 html_nodes(parsed_doc, xpath = "//div") %>% htmlAttrs()  
[[1]]  
      id          lang         date  
"R Inventor"    "english"   "June/2003"
```

[[2]]

lang	date
"english"	"October/2011"

end

Extraction of node elements

Attribute values

R code —

```
15 html_nodes(parsed_doc, xpath = "//div") %>% html_attr("lang")
[1] "english" "english"
```

— end

A silver lining

Do I really have to construct XPath expressions all by my own?

Do I really have to construct XPath expressions all by my own?

No!



Do I really have to construct XPath expressions all by my own?

No!



XPath creator tools

- *Selectorgadget*: <http://selectorgadget.com/>. Browser plugin that constructs XPath statements via a point-and-click approach. The generated expressions are not always efficient though
- *Web Developer Tools*: internal browser functionality which return XPath statements for selected nodes
- you will learn how to use these tools in upcoming sessions

Cologne Center for Comparative Politics

A Primer to Web Scraping with R

The Scraping Workflow

Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

Where you stand now

You have learned the main tools necessary to scrape static webpages with R

You have learned the main tools necessary to scrape static webpages with R

1. you are able to inspect HTML pages in your browser using the web developer tools

You have learned the main tools necessary to scrape static webpages with R

1. you are able to inspect HTML pages in your browser using the web developer tools
2. you are able to parse HTML into R with **rvest**

You have learned the main tools necessary to scrape static webpages with R

1. you are able to inspect HTML pages in your browser using the web developer tools
2. you are able to parse HTML into R with `rvest`
3. you are able to speak XPath

You have learned the main tools necessary to scrape static webpages with R

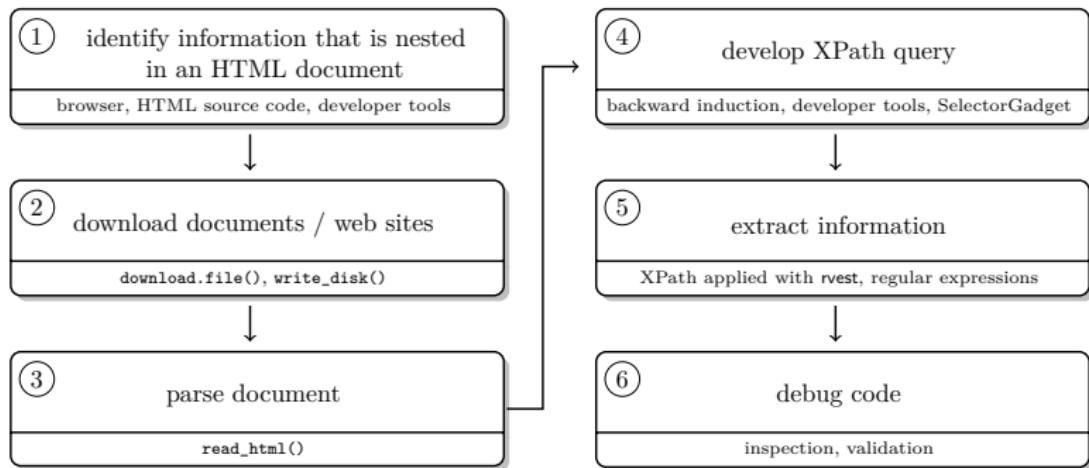
1. you are able to inspect HTML pages in your browser using the web developer tools
2. you are able to parse HTML into R with `rvest`
3. you are able to speak XPath
4. you are able to apply XPath expressions with `rvest`

You have learned the main tools necessary to scrape static webpages with R

1. you are able to inspect HTML pages in your browser using the web developer tools
2. you are able to parse HTML into R with `rvest`
3. you are able to speak XPath
4. you are able to apply XPath expressions with `rvest`
5. you are able to tidy web data with your R skills and regular expressions

The scraping workflow

The scraping workflow



Stay modest when accessing lots of data

- content on the web is publicly available, ...
- but accessing the data causes server traffic
- stay polite by querying resources as sparsely as possible

Downloading HTML files

Stay modest when accessing lots of data

- content on the web is publicly available, ...
- but accessing the data causes server traffic
- stay polite by querying resources as sparsely as possible

Two easy-to-implement practices

1. do not bombard the server with requests—and if you have to, do at a reasonable speed
2. download HTML files first, then parse

Downloading HTML files

R code

```
1 for (i in 1:length(list_of_urls)) {  
2     if (!file.exists(paste0(folder, file_names[i]))) {  
3         download.file(list_of_urls[i], destfile = paste0(folder, file_  
names[i]))  
4         Sys.sleep(runif(1, 1, 2))  
5     }  
6 }
```

end

Looping over a list of URLs

- `!file.exists()` checks whether a file does not yet exist in the local folder
- `download.file()` downloads the file to a folder; file name has to be specified
- `Sys.sleep()` suspends the execution of R code for a given time interval. Here: random interval between 1 and 2 seconds

Don't be a phantom

- downloading massive amounts of data may arouse attention from server administrators
- assuming that you've got nothing to hide, you should stay identifiable beyond your IP address

Don't be a phantom

- downloading massive amounts of data may arouse attention from server administrators
- assuming that you've got nothing to hide, you should stay identifiable beyond your IP address

Two easy-to-implement practices

1. personally get in touch with website owners
2. use HTTP header fields `From` and `User-Agent` to provide information about yourself

Staying identifiable

R code —

```
7 url <- "http://a-totally-random-website.com"
8 session <- html_session(url, add_headers(From = "my@email.com", `User-
Agent` = R.Version()$version.string)))
9 headlines <- session %>% html_nodes(xpath = "p//a") %>% html_text()

```

end

The code snippet explained

- `rvest`'s `html_session()` creates a session object that responds to HTTP and HTML methods
- here, we provide our email address and the current R version as User-Agent information
- this will pop up in the server logs—the webpage administrator has the chance to easily get in touch with you

Summary

Summary

- the basic scraping workflow with R is straightforward
- with great power comes great responsibility: stay polite on the web when scraping lots of data!
- more complexity is added when you want to gather data from multiple websites, or when dynamic elements such as forms or JavaScript content is involved
- we will consider such cases in upcoming sessions



Cologne Center for Comparative Politics

A Primer to Web Scraping with R

Dynamic Webpages

Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

Static webpages

The simple world of static webpages

Static webpages

- every user who visits the site gets the same content (unless the developer edits the source code)
- example: <https://www.jstatsoft.org>
- can contain "dynamic features", such as video or sound, but the page itself is not dynamic



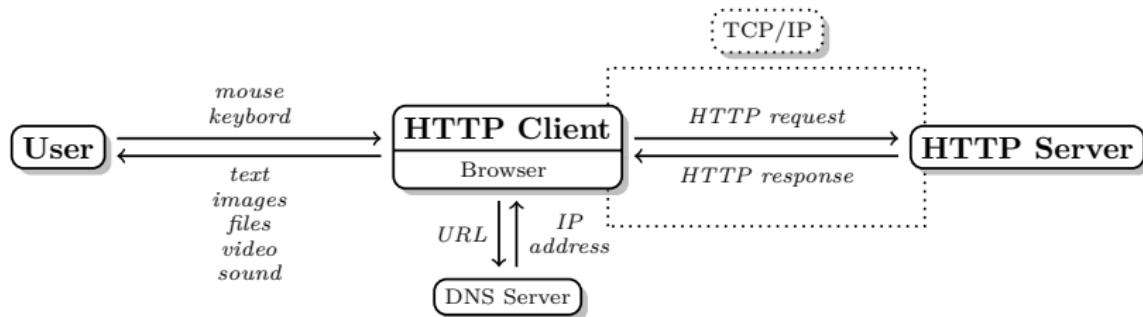
The screenshot shows the Amazon.com homepage from 1997. At the top, it says "WELCOME TO EARTH'S BIGGEST BOOKSTORE" and "Amazon.com". Below that, it displays "1.5 Million Books in Print" and "1 Million Out-of-Print Books". A large "A" logo is on the left. The main content area features a book cover for "The Great American Novel" by Hunter S. Thompson. To the right, there's a "First-Time Visitors Please Click Here" button. On the far right, there's a sidebar for "Book of the Day" featuring "The Great American Novel" again.

Screenshot of Amazon.com in 1997

A tiny bit of HTTP

Client-server communication with HTTP

- HTTP, the **Hypertext Transfer Protocol**, is a stateless protocol
- no information is retained by either sender or receiver
- makes interaction with websites straightforward, but not very exciting



Client-server communication with HTTP

1. Establishing connection

```
1 About to connect() to www.r-datacollection.com port 80 (#0)
2 Trying 173.236.186.125... connected
3 Connected to www.r-datacollection.com (173.236.186.125) port 80 (#0)
4 Connection #0 to host www.r-datacollection.com left intact
```

2. HTTP request

```
1 GET /index.html HTTP/1.1
2 Host: www.r-datacollection.com
3 Accept: */*
```

A tiny bit of HTTP

Client-server communication with HTTP

3. HTTP response

```
1  HTTP/1.1 200 OK
2  Date: Thu, 27 Feb 2014 09:40:35 GMT
3  Server: Apache
4  Vary: Accept-Encoding
5  Content-Length: 131
6  ...
7
8  <!DOCTYPE HTML PUBLIC "-//IETF//DTD HTML//EN">
9  <html> <head>
10 <title></title>
11 </head>
12 ...
```

4. Closing connection

```
1  Closing connection #0
```

The "problem" of static webpages

Static webpages reconsidered

- HTML/HTTP are used for static display of content → same content for every visitor
- in order to display dynamic content, they lack
 1. mechanisms to detect user behavior in the browser (and not only on the server)
 2. a scripting engine that reacts on this behavior
 3. a mechanism for asynchronous queries

Dynamic webpages

The not so simple world of dynamic webpages

Dynamic webpages

- with dynamic webpages, the displayed content can differ between users, even if the source code is the same
- things that can cause the display of different content:
 - operating system, browser, device
 - user actions on the page (mouse movements, scrolling, clicks, keyboard strokes)
 - conditions on the server-side



Source: <https://xkcd.com/869/> (Randall Munroe)

The not so simple world of dynamic webpages

Dynamic webpages

- massively used in modern webpage design and architecture
- (are thought to) enhance the user experience
- allow for much more ways to interact with webpage content

The not so simple world of dynamic webpages

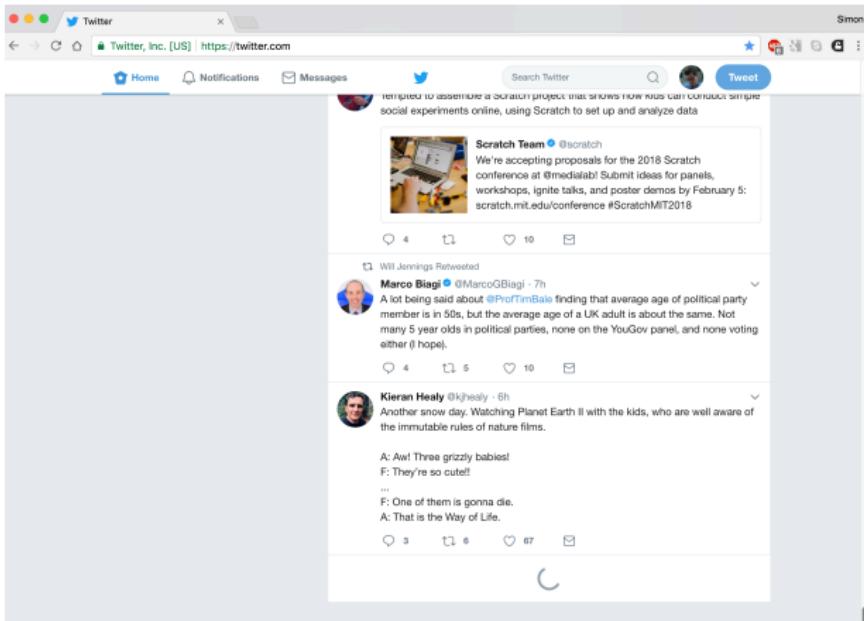
Dynamic webpages

- massively used in modern webpage design and architecture
- (are thought to) enhance the user experience
- allow for much more ways to interact with webpage content

Technologies

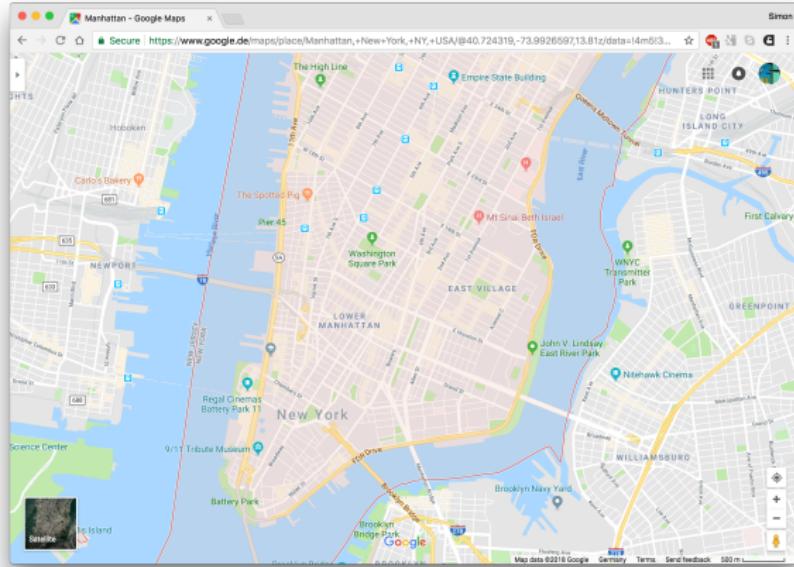
- in the case of **server-side** dynamic webpages, scripts on the server control webpage construction (e.g., PHP script reacts to user input to a form and returns subsets of a database)
- in the case of **client-side** dynamic webpages, (typically) JavaScript embedded in HTML determine what is displayed and how the DOM is changed
- a combination of these technologies is **Asynchronous JavaScript and XML (AJAX)**

Examples of dynamic webpages



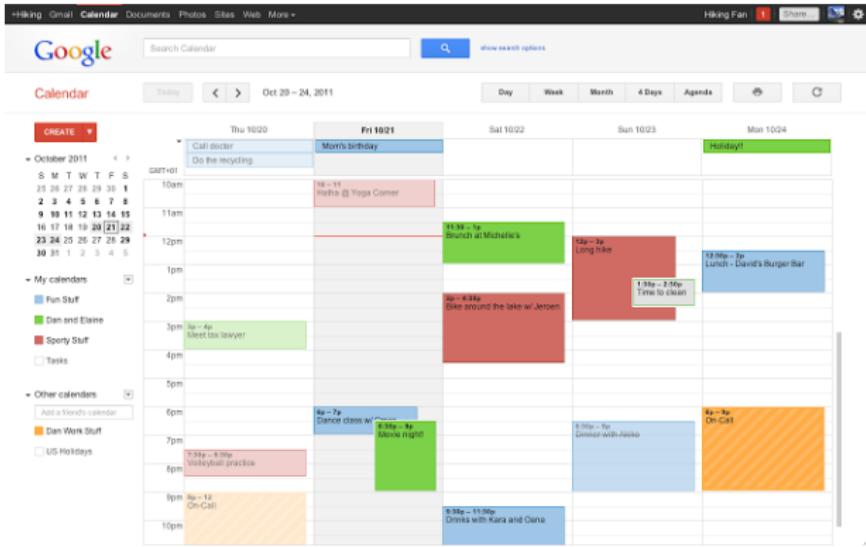
Your Twitter or Facebook feed that automatically gets updated when you scroll down

Examples of dynamic webpages



The map application that automatically loads new content when you zoom in

Examples of dynamic webpages



The calendar app that displays personal information and lets you interact a lot

The problem of dynamic webpages

The problem of dynamic webpages

The problem of dynamic webpages

- the tools you have encountered so far operated on the static source code: you access a page, download it, parse it

The problem of dynamic webpages

- the tools you have encountered so far operated on the static source code: you access a page, download it, parse it
- but what if content on the page changes, but the source code does not?

The problem of dynamic webpages

- the tools you have encountered so far operated on the static source code: you access a page, download it, parse it
- but what if content on the page changes, but the source code does not?
- dynamic webpages make classical screen scraping more difficult if not impossible

The problem of dynamic webpages

- the tools you have encountered so far operated on the static source code: you access a page, download it, parse it
- but what if content on the page changes, but the source code does not?
- dynamic webpages make classical screen scraping more difficult if not impossible
- and: dynamically rendered webpages become more and more common

The problem of dynamic webpages

- the tools you have encountered so far operated on the static source code: you access a page, download it, parse it
- but what if content on the page changes, but the source code does not?
- dynamic webpages make classical screen scraping more difficult if not impossible
- and: dynamically rendered webpages become more and more common
- before we learn about tools that can help us extract data in such scenarios, we will briefly consider the underlying processes, in particular AJAX technologies

Cologne Center for Comparative Politics

A Primer to Web Scraping with R

AJAX Technologies

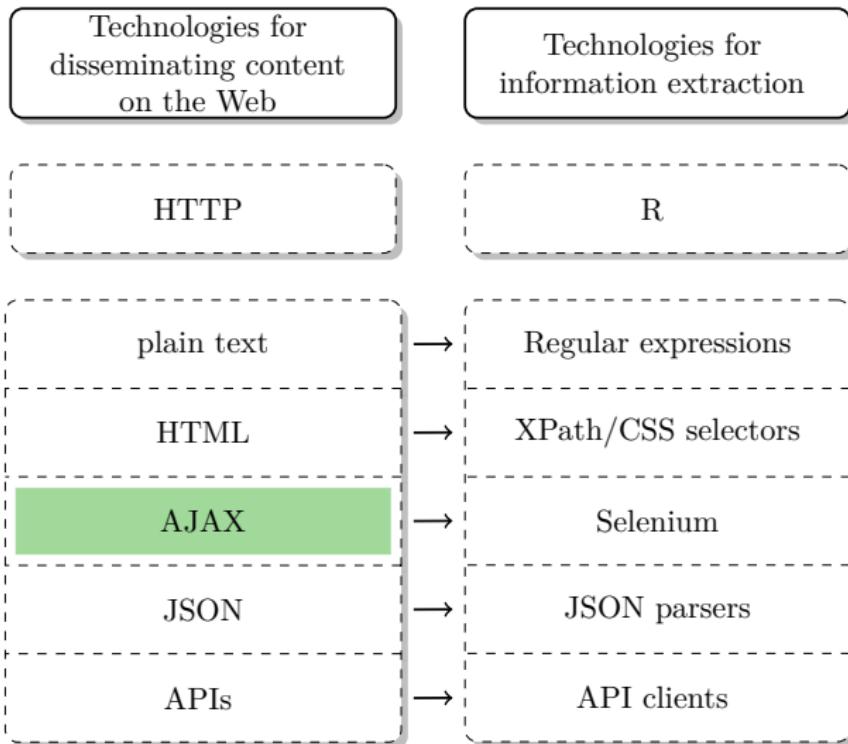
Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

AJAX

Technologies of the World Wide Web



Purpose

- recall: HTML/HTTP lack
 - 1. mechanisms to detect user behavior in the browser
 - 2. a scripting engine that reacts on this behavior
 - 3. a mechanism for asynchronous queries
- **A**synchronous **J**avaScript **a**nd **X**ML is a set of technologies that serve these purposes

What's AJAX?

Purpose

- recall: HTML/HTTP lack
 1. mechanisms to detect user behavior in the browser
 2. a scripting engine that reacts on this behavior
 3. a mechanism for asynchronous queries
- **A**synchronous **J**ava**S**cript **a**nd **X**ML is a set of technologies that serve these purposes

Components

- HTML/CSS for presentation
- Document Object Model (DOM; “tree structure”) to interact with data
- JSON/XML for data interchange
- XMLHttpRequest for asynchronous communication
- JavaScript as a scripting language

JavaScript

What's JavaScript? I

- Programming language that connects well to web technologies
- W3C web standard
- native browser support (built-in JS engine)
- nowadays employed on the majority of websites
- extensible by libraries, e.g. *jQuery*



What's JavaScript? I

- Programming language that connects well to web technologies
- W3C web standard
- native browser support (built-in JS engine)
- nowadays employed on the majority of websites
- extensible by libraries, e.g. *jQuery*



What's JavaScript? II

- technology to make webpages interactive ('dynamic HTML')
- allows to...
 - animate page elements
 - offer interactive content (games, video, ...)
 - manipulate page content and communication with the server without reloading the page

How's JavaScript code embedded in HTML?

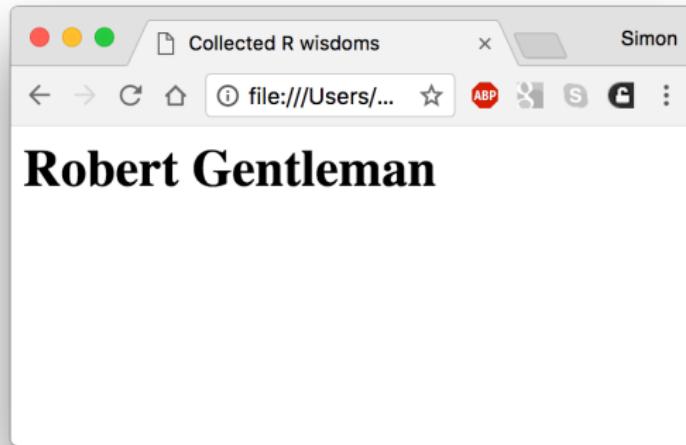
- between `<script>` tags
- as an external reference in the `src` attribute of a `<script>` element
- directly in certain HTML attributes ('event handler')

DOM manipulation with JavaScript

- adding/removing HTML elements
- changing attributes
- modification of CSS styles
- ...

Example:

```
1 <script type="text/javascript" src="jquery-1.8.0.min.js"></script>
2 <script type="text/javascript" src="script1.js"></script>
```



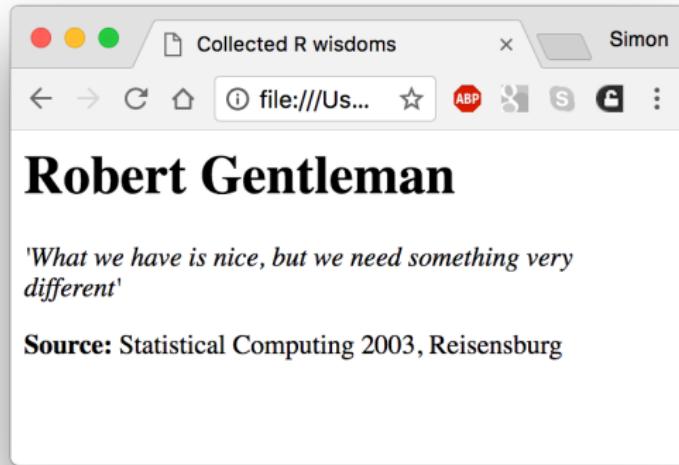
JavaScript on the Web

```
1 <!DOCTYPE HTML PUBLIC "-//IETF//DTD HTML//EN">
2 <html>
3
4 <script type="text/javascript" src="jquery-1.8.0.min.js"></script>
5 <script type="text/javascript" src="script1.js"></script>
6
7 <head>
8 <title>Collected R wisdoms</title>
9 </head>
10
11 <body>
12 <div id="R Inventor" lang="english" date="June/2003">
13   <h1>Robert Gentleman</h1>
14   <p><i>'What we have is nice, but we need something very different'</i>
15     </p>
16   <p><b>Source: </b>Statistical Computing 2003, Reisenburg</p>
17 </div>
18 </body>
19 </html>
```

A JavaScript code snippet

```
1 $(document).ready(function() {  
2     $("p").hide();  
3     $("h1").click(function(){  
4         $(this).nextAll().slideToggle(300);  
5     });  
6 });
```

- `$()` operator: addresses DOM elements
- `ready()`: JavaScript execution starts when the complete DOM is ready, i.e. fetched from the server
- `hide()`: element is hidden at first place
- `click` event: identifies mouse click and executes a certain action
- `nextAll()`: all subsequent elements in the DOM are addressed
- `slideToggle()`: Toggle effect, 300 milli-seconds



Summary

Summary

- AJAX technologies allow for a dynamic manipulation of the HTML tree
- what the user sees in the browser is not necessarily what's in the static HTML source code
- elements can be added, removed, or changed
- the "live" DOM tree that you saw in the Web Developer Tools takes track of such changes
- but you cannot simply access this live HTML with R and **rvest**
- we need another technology that helps us mimic a session in the browser to render all dynamic content



Cologne Center for Comparative Politics

A Primer to Web Scraping with R

Selenium: Basics

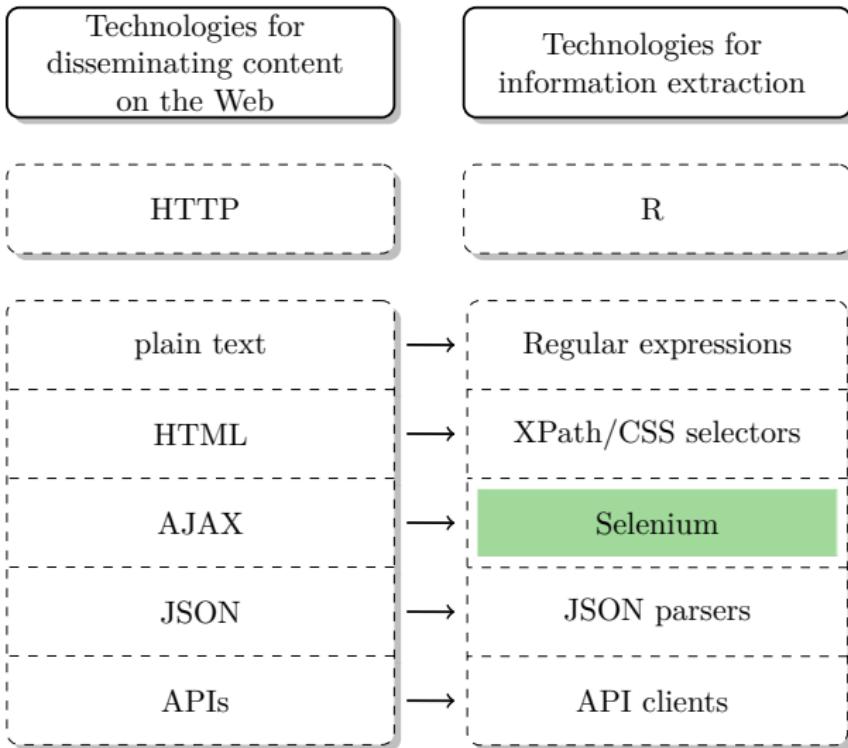
Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

Selenium

Technologies of the World Wide Web



The problem reconsidered

- dynamic data requests are not stored in the static HTML page
- therefore, we cannot access them with classical methods and packages (`httr`, `rvest`, `download.file()`, etc.)
- R is not a browser with a JavaScript rendering engine

Scraping dynamic webpages

The problem reconsidered

- dynamic data requests are not stored in the static HTML page
- therefore, we cannot access them with classical methods and packages (`httr`, `rvest`, `download.file()`, etc.)
- R is not a browser with a JavaScript rendering engine

The solution

- initiate and control a web browser session with R
- let the browser do the JavaScript interpretation work and the manipulations in the live DOM tree
- access information from the web browser session

What's Selenium?

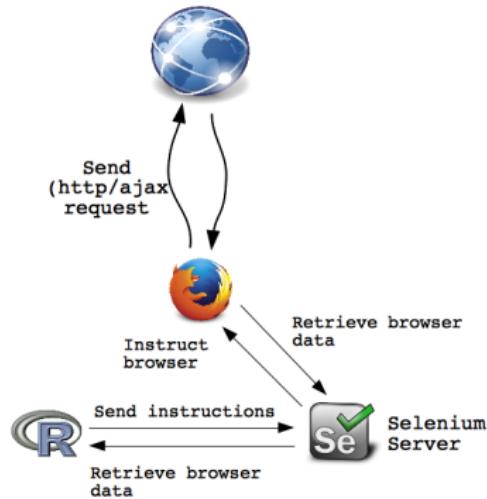
- free software environment for automated web application testing
- webpage: <http://www.seleniumhq.org>
- several modules for different tasks; most important for our purposes: Selenium WebDriver
- enables automated browsing via scripts



Selenium and R

The scraping workflow with Selenium and R

1. Selenium—an external, Java-based program—is launched
2. via the R package **RSelenium**, we remote-control Selenium and let it start a virtual server
3. we let Selenium start a browser session (e.g., Chrome)
4. everything we do in the browser—open a page, click on links or buttons, enter information—is described in R and sent to the server via the “remote-control” Selenium
5. we can gather the live HTML tree at any instance



Software requirements

- Java, <https://www.java.com/de/download/>
- Selenium Standalone Server, newest version available here:
<http://www.seleniumhq.org/download/> or via RSelenium and
`rsDriver()`
- browser (on most systems, both Chrome and Firefox seem to work reliably)
- **RSelenium** package

Cologne Center for Comparative Politics

A Primer to Web Scraping with R

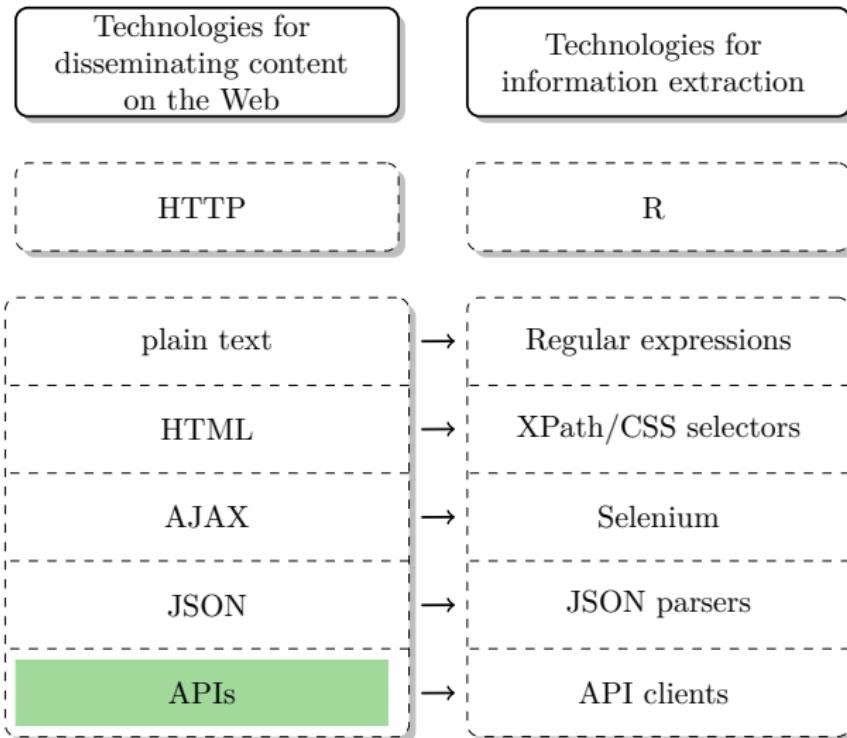
APIs

Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

Technologies of the World Wide Web



What are APIs?

What are APIs?

Definition

- **A**pplication **P**rogramming **I**nterface
- "data search engine": you pose a request, the API answers with a bulk of data
- let you/your program query a provider for specific data
- common data formats: XML, JSON
- many popular web services provide APIs (Twitter, Google, Facebook, Wikipedia, ...)

What are APIs?

Definition

- Application Programming Interface
- "data search engine": you pose a request, the API answers with a bulk of data
- let you/your program query a provider for specific data
- common data formats: XML, JSON
- many popular web services provide APIs (Twitter, Google, Facebook, Wikipedia, ...)

Why we should care about APIs

- provide instant access to clean data
- free us from building manual scrapers
- API usage implies mutual data collection agreement

Example

Example

Google Maps API

- Google provides access to powerful location services
- free service (at least when used modestly)
- input/output: places, names, coordinates, maps, ...
- see also: <https://developers.google.com/maps/documentation/>



Google Maps

Example

Google Maps API

- Google provides access to powerful location services
- free service (at least when used modestly)
- input/output: places, names, coordinates, maps, ...
- see also: <https://developers.google.com/maps/documentation/>



Google Maps

Potential use cases

- geocode observations based on address, city, post code, ...
- calculate distances between observations
- map observations

Access the API with R

- the `ggmap` package provides high-level functions to access API
- in this case, all we have to do is to figure out how the R function works – the communication with the API is processed completely in the background

Example

Access the API with R

- the **ggmap** package provides high-level functions to access API
- in this case, all we have to do is to figure out how the R function works – the communication with the API is processed completely in the background

R code

```
3 library(ggmap)
4 geocode("Berlin, Germany")
```

```
Information from URL : http://maps.googleapis.com/maps/api/geocode/json?
address=Berlin,%20Germany&sensor=false
      lon      lat
1 13.40495 52.52001
```

end

Example

Excerpt from raw JSON behind the call

```
1  {
2      "results" : [
3          {
4              "address_components" : [
5                  {
6                      "long_name" : "Berlin",
7                      "short_name" : "Berlin",
8                      "types" : [ "locality", "political" ]
9                  },
10                 ],
11                 "formatted_address" : "Berlin, Germany",
12                 "location" : {
13                     "lat" : 52.52000659999999,
14                     "lng" : 13.404954
15                 }
16             },
17             "place_id" : "ChIJAVkDPzd0qEcRcDteWOYgIQQQ",
18             "types" : [ "locality", "political" ]
19         }
20     ]
21 }
```

Example

Map the location

R code _____

```
5  get_googlemap("Berlin, Germany",
   zoom = 12, maptype = "hybrid")
%>% ggmap()
_____ end
```

Example

Map the location

R code

```
6 get_googlemap("Berlin, Germany",
  zoom = 12, maptype = "hybrid")
%>% ggmap()
----- end
```



Summary

Advantages of API use

- collecting data from the web using APIs provided by the data owner represents **the gold standard of web data retrieval**
- pure data collection without ‘layout waste’
- standardized data access
- de facto automatic agreement
- robustness of calls



[http://maxpixel.
freegreatpicture.com/
photo-1461569](http://maxpixel.freegreatpicture.com/photo-1461569)

Summary

Advantages of API use

- collecting data from the web using APIs provided by the data owner represents **the gold standard of web data retrieval**
- pure data collection without ‘layout waste’
- standardized data access
- de facto automatic agreement
- robustness of calls



[http://maxpixel.
freegreatpicture.com/
photo-1461569](http://maxpixel.freegreatpicture.com/photo-1461569)

Disadvantages of API use

- (sometimes) requires knowledge of API architecture
- dependent upon API suppliers
- use not always free

Cologne Center for Comparative Politics

A Primer to Web Scraping with R

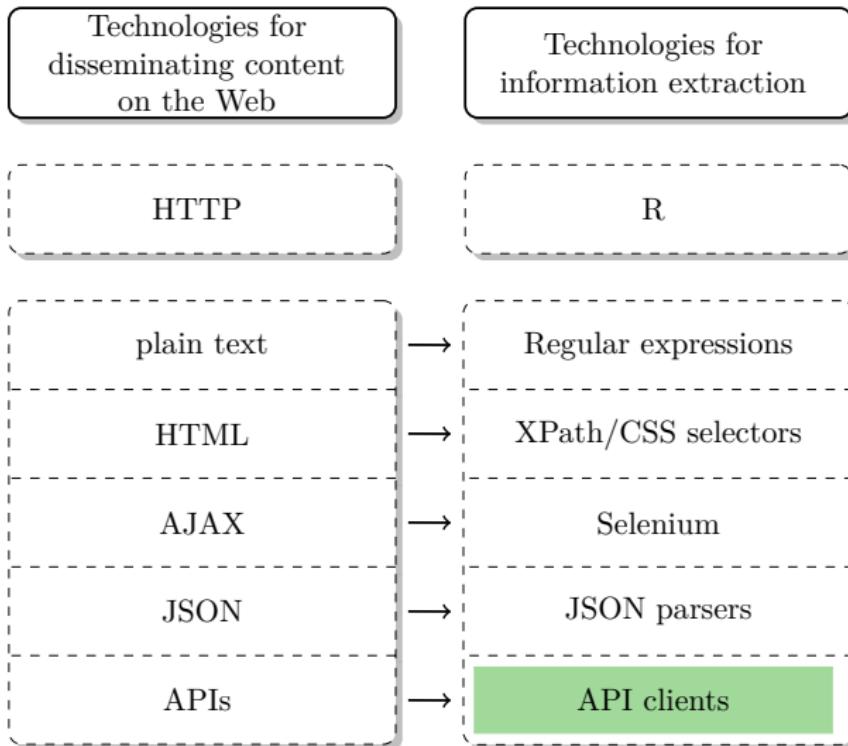
API Clients

Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

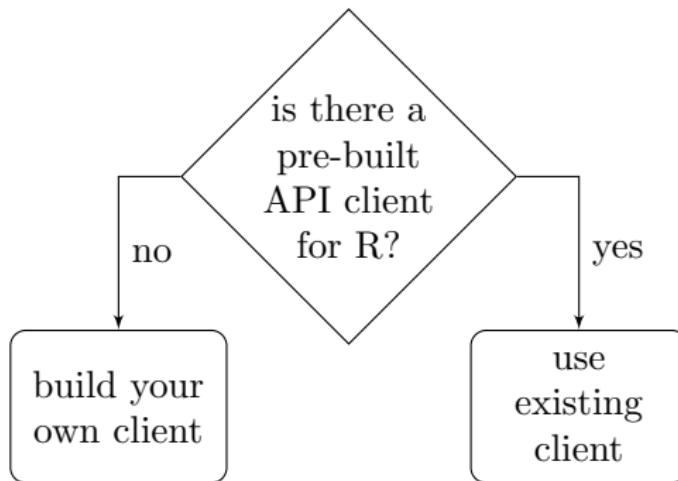
Technologies of the World Wide Web



Accessing APIs with R

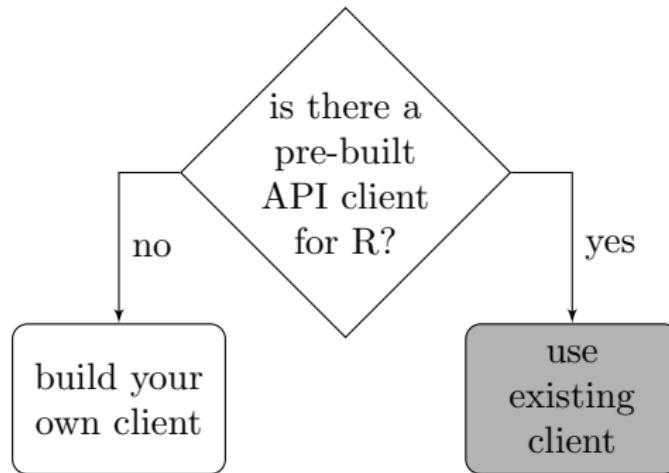
API access with R

There are basically **two scenarios:**



API access with R

There are basically **two scenarios:**



Here, we talk about the case where a client already exists.

API clients

- provide interface to APIs
- hide API back-end
- let you stay in your programming environment

API clients

- provide interface to APIs
- hide API back-end
- let you stay in your programming environment

Example

- the `rtweet` package provides an R client for Twitter
- lets you query data from the Twitter API with R commands
- you don't have to work with unfamiliar data formats (JSON) that is provided by the API—the client automatically transforms the incoming data into R objects

General resources

List of APIs: <http://www.programmableweb.com/apis>

rOpenSci – collection of R API clients:

<https://github.com/ropensci/opendata>

CRAN Task View:

<http://cran.r-project.org/web/views/WebTechnologies.html>

Finding API clients on the Web

General resources

List of APIs: <http://www.programmableweb.com/apis>

rOpenSci – collection of R API clients:

<https://github.com/ropensci/opendata>

CRAN Task View:

<http://cran.r-project.org/web/views/WebTechnologies.html>

How to find the API client you need

- google "R package + name of website"
- search on the website for a "Developer" or "API" section (only if you don't find an R package that works)

Popular R API clients

package name	access to	more info
<code>rtweet</code>	Twitter Stream and REST API	http://rtweet.info/
<code>Rfacebook</code>	Facebook API	https://github.com/pablobarbera/Rfacebook
<code>ipapi</code>	ip-api.com's API	https://github.com/hrbrmstr/ipapi
<code>ggmap</code>	Google Maps/OpenStreetMap APIs	https://github.com/dkahle/ggmap
<code>eurostat</code>	Eurostat database	https://github.com/ropengov/eurostat
<code>rftimes</code>	New York Times APIs	https://cran.rstudio.com/web/packages/rftimes/index.html

Summary

Summary

- if the website you want to use as a data basis for your research provides access via an API, you can consider yourself lucky
- if there is a working R-based API client available, you hit the jackpot



Summary

- if the website you want to use as a data basis for your research provides access via an API, you can consider yourself lucky
- if there is a working R-based API client available, you hit the jackpot



Final considerations

- please always cite package authors when you use their work. run `citation("<package name>")` to see how
- sometimes API clients are not up to date. Consider to notify the author or help update the package
- sometimes API clients provide functionality for only a fraction of the API's capability. It's always worth to check out the API documentation if you are looking for specific features

Cologne Center for Comparative Politics

A Primer to Web Scraping with R

Accessing APIs from Scratch

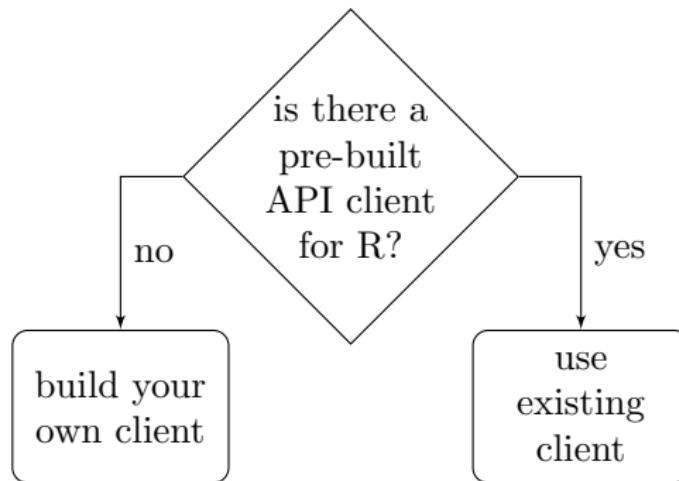
Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

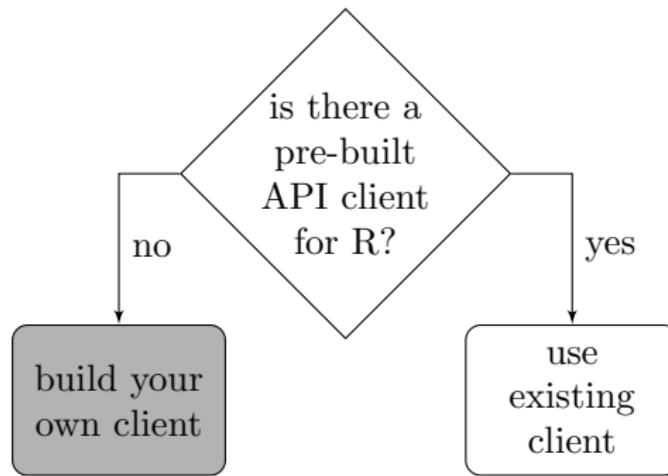
API access with R

There are basically **two scenarios**:



API access with R

There are basically **two scenarios:**



Here, we talk about the case where you have to build your own API binding.

Accessing APIs from Scratch

Accessing APIs from Scratch

Steps to be taken

Steps to be taken

1. figure out how the API works: every API has a human-readable documentation for developers

Steps to be taken

1. figure out how the API works: every API has a human-readable documentation for developers
2. build access to the API via R

Steps to be taken

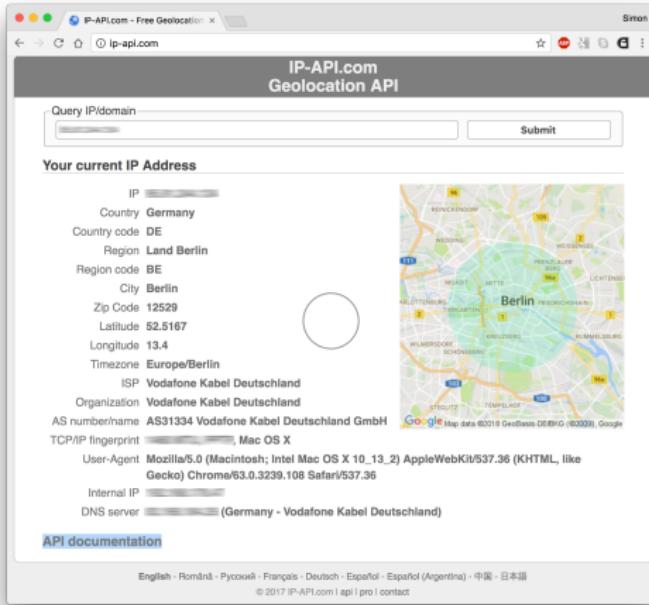
1. figure out how the API works: every API has a human-readable documentation for developers
2. build access to the API via R
3. build functionality that processes the data output (e.g., turns it into R objects)

Example

Example

The IP API

- available at
<http://ip-api.com/>
- its purpose: parses your IP (or a bunch of given IPs) and returns some values, including guessed location, service provider, and organization behind the IP



The screenshot shows the IP-API.com geolocation API interface. At the top, there's a search bar labeled "Query IP/domain" with the placeholder "125.29.52.167" and a "Submit" button. Below the search bar, it says "Your current IP Address". To the right is a map of Berlin, Germany, with various districts like Mitte, Kreuzberg, and Friedrichshain highlighted in green and yellow. A circular placeholder icon is positioned over the map. On the left side of the map, there's a list of geographical and technical details:

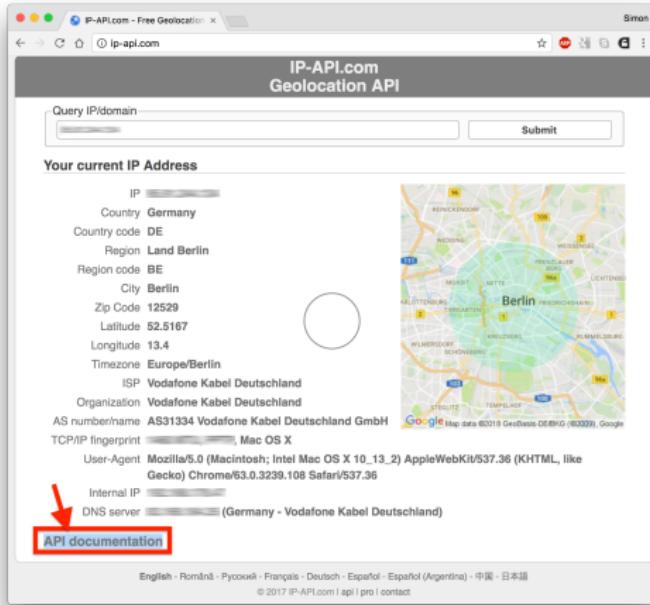
- IP 125.29.52.167
- Country Germany
- Country code DE
- Region Land Berlin
- Region code BE
- City Berlin
- Zip Code 12529
- Latitude 52.5167
- Longitude 13.4
- Timezone Europe/Berlin
- ISP Vodafone Kabel Deutschland
- Organization Vodafone Kabel Deutschland
- AS number/name AS31334 Vodafone Kabel Deutschland GmbH
- TCP/IP fingerprint [REDACTED] Mac OS X
- User-Agent Mozilla/5.0 (Macintosh; Intel Mac OS X 10_13_2) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/63.0.3239.108 Safari/537.36
- Internal IP [REDACTED]
- DNS server [REDACTED] (Germany - Vodafone Kabel Deutschland)

At the bottom of the page, there's a link to "API documentation" and a footer with language links: English · Română · Русский · Français · Deutsch · Español · Español (Argentina) · 中文 · 日本語. The footer also includes copyright information: © 2017 IP-API.com | api | pro | contact.

Example

The IP API

- available at
<http://ip-api.com/>
- its purpose: parses your IP (or a bunch of given IPs) and returns some values, including guessed location service provider, and organization behind the IP
- check out API documentation



The screenshot shows the IP-API.com geolocation API interface. A red arrow points to the "API documentation" link at the bottom of the page.

Query IP/domain: [REDACTED]

Your current IP Address

IP: [REDACTED]

Country: Germany

Country code: DE

Region: Land Berlin

Region code: BE

City: Berlin

Zip Code: 12529

Latitude: 52.5167

Longitude: 13.4

Timezone: Europe/Berlin

ISP: Vodafone Kabel Deutschland

Organization: Vodafone Kabel Deutschland

AS number/name: AS31334 Vodafone Kabel Deutschland GmbH

TCP/IP fingerprint: [REDACTED], Mac OS X

User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10_13_2) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/63.0.3239.108 Safari/537.36

Internal IP: [REDACTED]

DNS server: [REDACTED] (Germany - Vodafone Kabel Deutschland)

[API documentation](#)

English · Română · Русский · Français · Deutsch · Español · Español (Argentina) · 中国 · 日本語
© 2017 IP-API.com | [api](#) | [pro](#) | [contact](#)

Example

The IP API

- access is free for non-commercial use (up to 150 requests per minute)
- various response formats available
- paid pro service available

The screenshot shows a web browser window displaying the documentation for the IP Geolocation API at ip-api.com/docs/. The page title is "IP Geolocation API". On the left, there's a sidebar with a "Table of Contents" section containing links to "IP Geolocation API", "About", "Response formats and examples", and "Usage limits". The main content area has several sections: "About" (describing free usage), "Response formats and examples" (listing CSV, JSON, Newline Separated, Serialized PHP, XML, Change Log, Statistics, DNS API, IP Geolocation API, and Urban IP), "Usage limits" (warning about a ban for over 150 requests/min), and a note for commercial users. At the bottom right is a "Back to top" button.

Example

The IP API

- to access the API, we have to send a **GET** request to <http://ip-api.com/json>
- we can provide any IP we want
- the response is raw JSON code

The screenshot shows a web browser displaying the IP Geolocation API documentation at <http://ip-api.com/docs/json>. The page has a sidebar with links to CSV, Error messages, JSON, Batch JSON, Nested, Separated, Return values, Serialized PHP, XML, Change Log, Statistics, DNS API, IP Geolocation API, and Urban IP. The main content area is titled "JSON" and contains sections for "Usage" and "Response". Under "Usage", it says "To receive the response in JSON format, send a GET request to" and provides a link to <http://ip-api.com/json>. Under "Response", it says "A successful request will return, by default, the following:" and shows a JSON object structure. Below the JSON object, there is a "List of returned values" section and a note about failed requests.

```
{  
    "status": "success",  
    "country": "COUNTRY",  
    "countryCode": "COUNTRY CODE",  
    "region": "REGION CODE",  
    "regionName": "REGION NAME",  
    "city": "CITY",  
    "zip": "ZIP CODE",  
    "lat": LATITUDE,  
    "lon": LONGITUDE,  
    "timezone": "TIME ZONE",  
    "isp": "ISP NAME",  
    "org": "ORGANIZATION NAME",  
    "as": "AS NUMBER / NAME",  
    "query": "IP ADDRESS USED FOR QUERY"  
}
```

List of returned values

A failed request will return, by default, the following:

Example

IP API access in R

A **GET** request essentially means that we assemble a URL using

- the API **endpoint** (a basic URL) and
- **parameters-value pairs** that are added to the URL
- here, we only add the literal IP address
- we use **fromJSON** to pose the request and directly parse the data from the API

R code

```
1 url <- "http://ip-api.com/json"
2 url_ip <- paste0(url, "/208.80.152.201")
3 ip_parsed <- jsonlite::fromJSON(url_ip)
```

end

Example

Investigating the output

R code

```
4 names(ip_parsed)
[1] "as"          "city"        "country"      "countryCode"  "isp"
[6] "lat"         "lon"         "org"         "query"       "region"
[11] "regionName" "status"      "timezone"    "zip"

5 ip_parsed
$as
[1] "AS14907 Wikimedia Foundation Inc."

$city
[1] "San Francisco (South Beach)"

$country
[1] "United States"

$countryCode
[1] "US"
```

Example

Bringing it into shape

- in our case, `fromJSON()` has created an R list object
- we turn it into a data frame

R code

```
6 ip_parsed %>% as.data.frame(stringsAsFactors = FALSE)
               as                      city
1 AS14907 Wikimedia Foundation Inc. San Francisco (South Beach)
      country countryCode           isp     lat     lon
1 United States             US Wikimedia Foundation 37.787 -122.4
      org          query region regionName   status
1 Wikimedia Foundation 208.80.152.201       California success
      timezone    zip
1 America/Los_Angeles 94105
```

end

Summary

Summary

- accessing APIs from scratch can be almost as easy as using a readily available R client
- often however, APIs are more complex (provide much more parameters and variations in data output)
- ideally, you build functions around your API calls to make it even more accessible



Cologne Center for Comparative Politics

A Primer to Web Scraping with R

Legal Issues

Simon Munzert

Hertie School of Governance, Berlin

November 12, 2018

Is web scraping legal?

Is web scraping legal?

Disclaimer

Obviously, I am not a lawyer. Do not rely on any of my comments on this topic. If you are seriously worried about the legality of your scraping work, please consult a legal expert.



Is web scraping legal?

Scraping vs. crawling

- **Web scraping:** downloading data from a very specific page in a (semi-)automated manner
- **Web crawling:** automatically downloading webpage data, extracting hyperlinks, following links, downloading webpage data, ... (e.g.: Googlebot, BaiduSpider)

This course mostly is scraping, not crawling or web harvesting!

Is web scraping legal?

Scraping vs. crawling

- **Web scraping:** downloading data from a very specific page in a (semi-)automated manner
- **Web crawling:** automatically downloading webpage data, extracting hyperlinks, following links, downloading webpage data, ... (e.g.: Googlebot, BaiduSpider)

This course mostly is scraping, not crawling or web harvesting!

- web scraping per se is not illegal
- there is no unambiguous **yes** or **no** for specific applications in any country according to current jurisdiction
- so far, legal cases (especially in the US) often (but not always) dealt with commercial interest, crawling applications, and often (but not always) huge masses of data, e.g., *eBay vs. Bidder's Edge*, *AP vs. Meltwater*, *Facebook vs. Pete Warden*

Is web scraping legal?

Why web scraping can be problematic

- violation of copyrights, Terms of Service, consumption of bandwidth

Is web scraping legal?

Why web scraping can be problematic

- violation of copyrights, Terms of Service, consumption of bandwidth

Some counter-arguments

- data is publicly accessible
 - but the website as a "creative arrangement" might be copyrighted
- this is fair use → depends on your use
- it's the same what my browser does
 - not exactly, and many ToS prohibit automated uses of their data
- this is unfair—Google's business model is built on crawling the whole web
 - true, but you are not Google

Some interesting comments by Pablo Hoffman, co-founder of Scrapinghub

- As long as they don't crawl at a disruptive rate, scrapers do not breach any contract (in the form of terms of use) or commit a crime (as defined in the Computer Fraud and Abuse Act).
- Website's user agreement is not enforceable as a browse-wrap agreement because companies do not provide sufficient notice of the terms to site visitors.
- Scrapers accesses website data as a visitor, and by following paths similar to a search engine. This can be done without registering as a user (and explicitly accepting any terms).

(found on <https://goo.gl/jTt4ER>)

Things webpage administrators can do to prevent you from scraping massive amounts of data from their pages

- block your IP address
- identify your approximate geolocation from your IP address, then block
- move content exclusively to web services / APIs
- block bots with a particular user agent string (more on that later)
- challenge-response tests like **C**ompletely **A**utomated **P**ublic **T**uring test to tell **C**omputers and **H**umans **A**part (CAPTCHA)
- obfuscation of data
- frequent changes in HTML/CSS

Summary

Summary

- there is no unconditional "legal" or "illegal" status of web scraping
- your use of the data can violate the data owner's rights
- targeted scraping efforts with limited traffic usually do not cause any problems



Cologne Center for Comparative Politics

A Primer to Web Scraping with R

Good Practice

Simon Munzert

Hertie School of Governance, Berlin

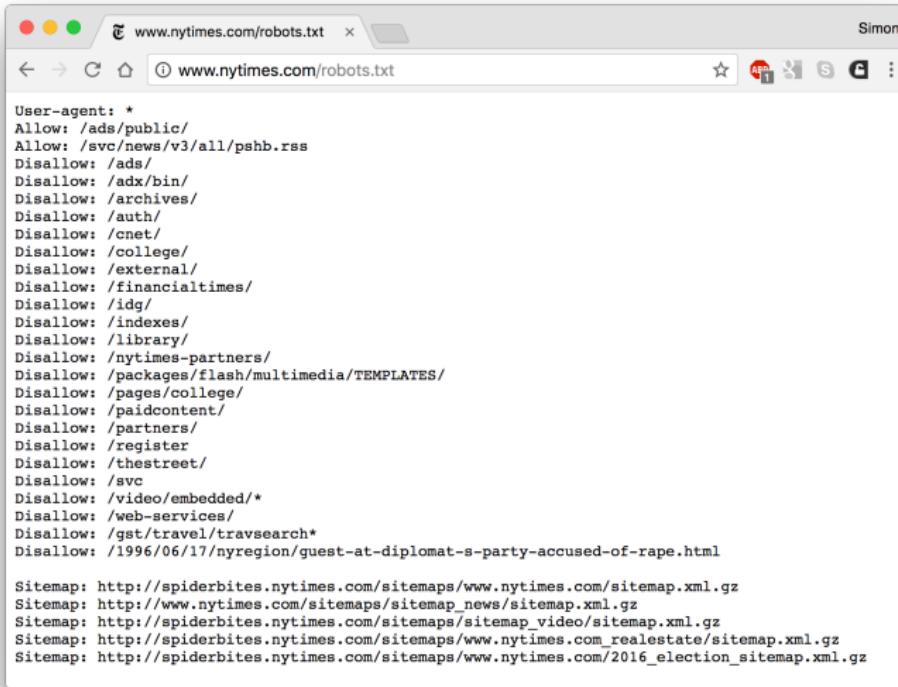
November 12, 2018

Understanding robots.txt

What's robots.txt?

- ‘Robots Exclusion Protocol’, informal protocol to prohibit web robots from crawling content
- located in the root directory of a website (e.g.,
<http://www.google.com/robots.txt>)
- documents which bot is allowed to crawl which resources (and which not)
- not a technical barrier, but a sign that asks for compliance

Example: NY Times's robots.txt



A screenshot of a web browser window titled "www.nytimes.com/robots.txt". The address bar also shows "www.nytimes.com/robots.txt". The page content displays the NY Times's robots.txt file. The file contains several "User-agent: *" entries with various "Allow:" and "Disallow:" directives. It also includes sections for "Sitemap" URLs.

```
User-agent: *
Allow: /ads/public/
Allow: /svc/news/v3/all/pshb.rss
Disallow: /ads/
Disallow: /adx/bin/
Disallow: /archives/
Disallow: /auth/
Disallow: /cnet/
Disallow: /college/
Disallow: /external/
Disallow: /financialtimes/
Disallow: /idg/
Disallow: /indexes/
Disallow: /library/
Disallow: /nytimes-partners/
Disallow: /packages/flash/multimedia/TEMPLATES/
Disallow: /pages/college/
Disallow: /paidcontent/
Disallow: /partners/
Disallow: /register/
Disallow: /thestreet/
Disallow: /svc
Disallow: /video/embedded/*
Disallow: /web-services/
Disallow: /gst/travel/travsearch*
Disallow: /1996/06/17/nyregion/guest-at-diplomat-s-party-accused-of-rape.html

Sitemap: http://spiderbites.nytimes.com/sitemaps/www.nytimes.com/sitemap.xml.gz
Sitemap: http://www.nytimes.com/sitemaps/sitemap_news/sitemap.xml.gz
Sitemap: http://spiderbites.nytimes.com/sitemaps/sitemap_video/sitemap.xml.gz
Sitemap: http://spiderbites.nytimes.com/sitemaps/www.nytimes.com_realestate/sitemap.xml.gz
Sitemap: http://spiderbites.nytimes.com/sitemaps/www.nytimes.com/2016_election_sitemap.xml.gz
```

Syntax in robots.txt

Syntax

- not an official W3C standard, partly inconsistent syntax
- rules listed bot by bot
- general, bot-independent rules under '*' (most interesting entry for R-based crawlers)
- directories/folders listed separately

```
1 User-agent: Googlebot
2 Disallow: /images/
3 Disallow: /private/
```

```
1 User-agent: *
2 Disallow: /private/
```

Syntax in robots.txt

Universal ban

```
1 User-agent: *
2 Disallow: /
```

Separation of bots by empty line

```
1 User-agent: Googlebot
2 Disallow: /images/
4 User-agent: Slurp
5 Disallow: /images/
```

Allow declaration

```
1 User-agent: *
2 Disallow: /images/
3 Allow: /images/public/
```

Syntax in robots.txt

Crawl-delay (in seconds)

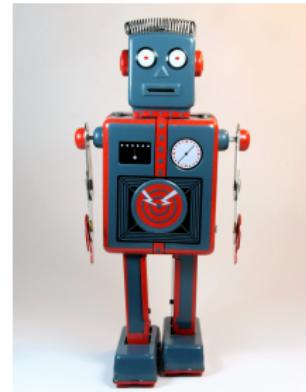
```
1 User-agent: *
2 Crawl-delay: 2
3 User-Agent: Googlebot
4 Disallow: /search/
```

Robots <meta> tag

```
1 <meta name="robots" content="noindex,nofollow" />
```

How to deal with robots.txt?

- not clear if robots.txt is legally binding or not, and if yes for which activities
- originally not thought of as protection against small-scale web scraping applications, but against large-scale indexing bots
- guide to a webmaster's preferences with regards to visibility of content
- my advice: take `robots.txt` into account! If the data you are interested in are excluded from crawling: contact webmaster
- there is even an R package, `robotstxt`, that helps you parse `robots.txt` documents and stick to the rules. It's available here:
<https://github.com/ropenscilabs/robotstxt>



By D J Shin - My Toy Museum, [https://commons.wikimedia.org/wiki/File:QSH_Tin_Wind_Up_Mechanical_Robot_\(Giant_Easelback_Robot\)_Front.jpg](https://commons.wikimedia.org/wiki/File:QSH_Tin_Wind_Up_Mechanical_Robot_(Giant_Easelback_Robot)_Front.jpg)

Scraping etiquette

Some advice for your work

1. You take all the responsibility for your web scraping work.

Some advice for your work

1. You take all the responsibility for your web scraping work.
2. Take all copyrights of a country's jurisdiction into account.

Some advice for your work

1. You take all the responsibility for your web scraping work.
2. Take all copyrights of a country's jurisdiction into account.
3. If you publish data, do not commit copyright fraud.

Some advice for your work

1. You take all the responsibility for your web scraping work.
2. Take all copyrights of a country's jurisdiction into account.
3. If you publish data, do not commit copyright fraud.
4. If possible, stay identifiable.

Some advice for your work

1. You take all the responsibility for your web scraping work.
2. Take all copyrights of a country's jurisdiction into account.
3. If you publish data, do not commit copyright fraud.
4. If possible, stay identifiable.
5. If in doubt, ask the author/creator/provider of data for permission—if your interest is entirely scientific, chances aren't bad that you get data.

Some advice for your work

1. You take all the responsibility for your web scraping work.
2. Take all copyrights of a country's jurisdiction into account.
3. If you publish data, do not commit copyright fraud.
4. If possible, stay identifiable.
5. If in doubt, ask the author/creator/provider of data for permission—if your interest is entirely scientific, chances aren't bad that you get data.
6. Consult current jurisdiction, e.g. on
<http://blawgsearch.justia.com> or from a lawyer specialized on internet law.

Scraping etiquette

