



RVest Package

By Gregory Albarian



Installation

```
install.packages("rvest")
```

```
library(rvest)
```

```
# used for web scraping but can also do some web  
scraping
```

Basic HTML Syntax

```
<!DOCTYPE html>  
<html>  
  <head>  
    <meta charset="UTF-8">  
    <title>blank site</title>  
    <!--comments-->  
  </head>  
  <body>  
    <p>"Paragraphs"</p>  
  </body>  
</html>
```



Example.com body example

```
</head>

<body>
  <div>
    <h1>Example Domain</h1>
    <p>This domain is for use in illustrative examples in documents. You may use this
    domain in literature without prior coordination or asking for permission.</p>
    <p><a href="https://www.iana.org/domains/example">More information...</a></p>
  </div>
</body>
</html>
```



Example Domain

This domain is for use in illustrative examples in documents. You may use this domain in literature without prior coordination or asking for permission.

[More information...](#)

Full-screen Snip



Scraping Example.com

```
html_code = read_html("http://www.example.com/")
```

%>% is like a pipe it says we are still referring to the same data in the next line

```
html_code %>%
```

```
html_nodes("p") %>% #use node to select one and nodes to select all
```

```
html_text()
```

#one side note you can find tables but putting into data frames and parsing them well needs other libraries

The form of a table - goes in the body

`<table style="width:50%">`

`<tr>`

`<th>1,1</th>`

`<th>1,2</th>`

`<th>1,3</th>`

`</tr>`

`<tr>`

`<td>2,1</td>`

`<td>2,2</td>`

`<td>2,3</td>`

`</tr>`

`<tr>`

`<td>3,1</td>`

`<td>3,2</td>`

`<td>3,3</td>`

`</tr>`

`</table>`

Display from last slide

	1,1	1,2	1,3
2,1			
3,1			



Different ways to scrape tables

Practical Example scraping the stock information for the DOW Jones on Yahoo! Finance

One way: - use the `html_node()` or `html_nodes()` functions like we have been doing

Second way: - use a special builtin function - can store as list in a variable

```
html_table(read_html("https://finance.yahoo.com/quote/%5EDJI?p=DJI&tsrc=fin-srch"))
```



Output

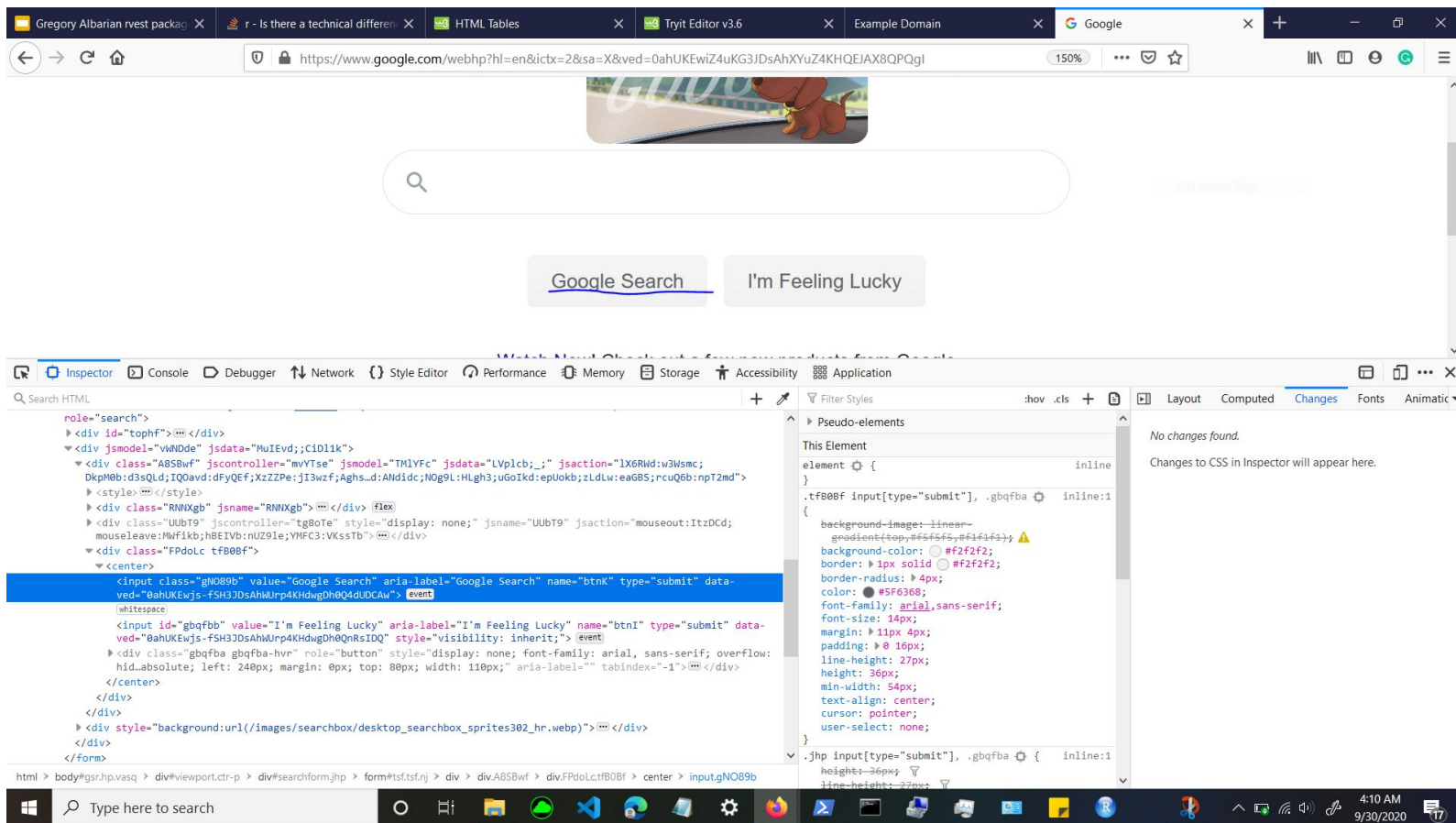
- The output is stored as a list
- The list is full of String the numbers need to be converted to numeric values
- Can be used for more practical data analysis



Scraping Forms - finding user input on sites

- Again you can just use the `get_nodes()` function
- There is a special function
- Later example <https://google.com/>

Practical example - Google Search



The screenshot shows a web browser window displaying the Google search page. The address bar shows the URL: `https://www.google.com/webhp?hl=en&icx=2&sa=X&ved=0ahUKEwiZ4uKG3JDsAhXYuZ4KHQJAX8PQqgl`. The page features the Google logo, a search bar, and two buttons: "Google Search" and "I'm Feeling Lucky".

The browser's developer tools are open, showing the HTML and CSS for the search button. The HTML structure is as follows:

```
<div class="gNO89b" value="Google Search" aria-label="Google Search" name="btnK" type="submit" data-ved="0ahUKEwjs-fSH3JDsAhMjrp4KHdwgDh0Q4dUDCAu"> event
  <input id="gqfbb" value="I'm Feeling Lucky" aria-label="I'm Feeling Lucky" name="btnI" type="submit" data-ved="0ahUKEwjs-fSH3JDsAhMjrp4KHdwgDh0Q4dUDCAu"> event
  <div class="gbqfba gbqfba-hvr" role="button" style="display: none; font-family: arial, sans-serif; overflow: hidden; absolute; left: 240px; margin: 0px; top: 80px; width: 110px;" aria-label="" tabindex="-1">
    </div>
  </div>
  <div style="background:url(/images/searchbox/desktop_searchbox_sprites302_hr.webp)">
    </div>
  </form>
```

The CSS for the search button is shown in the right pane:

```
.tfB08F input[type="submit"], .gbqfba {
  background-image: linear-gradient(to top, #f5f5f5, #f5f5f5);
  background-color: #f2f2f2;
  border: 1px solid #f2f2f2;
  border-radius: 4px;
  color: #5f6368;
  font-family: arial, sans-serif;
  font-size: 14px;
  margin: 11px 4px;
  padding: 0 16px;
  line-height: 27px;
  height: 36px;
  min-width: 54px;
  text-align: center;
  cursor: pointer;
  user-select: none;
}
```

The browser's taskbar at the bottom shows the Windows logo, a search bar, and several application icons. The system clock in the bottom right corner displays the time as 4:10 AM on 9/30/2020.



Code:

```
url <-  
"https://www.google.com/webhp?hl=en&sa=X&ved=0ahUKEwiAquTZ8IzsAhUOrJ4  
KHaesD7EQPAgl"  
  
html_form(read_html(url))  
  
# the output should show the different user inputs on the page  
  
# allows us to see the Google Search button in the code
```



Sources:

- https://rdr.io/cran/rvest/man/html_nodes.html#heading-3
- <https://github.com/tidyverse/rvest>
- <http://rvest.tidyverse.org/>
- <https://www.dataquest.io/blog/web-scraping-in-r-rvest/>
- <https://cran.r-project.org/web/packages/rvest/rvest.pdf>