# Workshop 4: Fundamentals of Web Scraping

## Myles D. Garvey, Ph.D

### Winter, 2020

## 1 Learning Objectives

1. Learn the fundamentals of HTML, CSS, and Javascript.

2. Learn how data is represented in an HTML DOM Tree

3. Learn how to "read" HTML code in R and how to extract information from HTML.

4. Learn the difference between static and dynamic web scrapping.

5. Learn how to use the `rvest` package to conduct static web scrapping.

6. Learn how to use the `RSelenium` package to conduct dynamic web scrapping.

## 2 Workshop Description and Submission

While APIs are very useful technologies to be able to obtain data, the truth of the matter is that most good data for which we need is out and about on websites, unavailable via API. In some instances, we have such a great need for data that is not in an API, we need to resort to manual collection of information. For example, suppose we would like to gather a list of contacts from every business school in the country. If we have a list of websites to visit, we could go one by one, visiting the website, going to the faculty directory, and copy and pasting the faculty's information into an Excel file.

Needless to say, this process would take us a very long time. However, with a little bit of knowledge on how websites work in the first place, one could look at common patterns in these websites, and download the information automatically. A software that accomplishes such a task is called a *web scrapper*, and it is a very useful tool that can save us hours if not weeks of work. Unfortunately, given that the needs of data, and where the data is stored, and the format in which it is stored, varies very greatly, it is almost impossible to have a single software out there that can accomplish our goal. Hence, we usually must resort to *writing our own web scrapper*. This usually entails a simple process of first understanding the structure of the website from which we would like to scrape (if we are only doing so from a single website), and second, interacting with the website to download the data we need.

In this workshop, you will learn how websites "work". You do not need to become expert website designers. But you will need to know (1) how information on a website is represented, (2) how information on a website is instructed to be displayed and (3) how information can change on a website due to a user's interaction with the site. Once you have a fundamental understanding of this, only then can you move to web scraping. The subsequent tutorials will teach you how to extract web information in R, as well as how to scrape the web in a static and dynamic manner using R.

You will submit to blackboard a .PDF document. Complete each exercise in this workshop. These are easy and small, and so it should not take you very long. In a word document, type up the question number, the original question text, and your response in R code. Save your word document as a .PDF and submit it to the submission link on Blackboard. Please note that failure to submit this as a .PDF will result in a 0 on the workshop.

## 3   Tutorial 1: Understanding HTML, CSS, DOM, and Javascript

In order to understand how to scrape information from a website, one must first understand how a website is structured in the first place. Websites (with the exception of Flash Websites, which are rarely a thing these days) are constructed using three major technologies: Hyper-Text Markup Language (HTML), Cascading Style Sheets (CSS), and Javascript. With these three technologies, a website is nothing more than (1) what elements should be displayed on the site, (2) how and where should they be displayed, and (3) how can the website change based on user interaction? When a creater of a website has specified the information to be shown, the way in which the information should be shown, and how the user can interact with the website to change the information, they do so by writing HTML/CSS/JS code. This code all feeds into a software called a browser. The browser dissects this information and displays all of it according to the creator's design. Therefore, for a basic website, the designer must:

- Specify the HTML code, which will tell the browser what specific information is to be displayed on the screen of the browser.

- Specify the CSS code, which will tell the browser what colors to use, sizes of boxes, pictures and media, as well as menus and general "look and feel" of the website. Put simply, the CSS is code that tells the browser *how* the information specified in the HTML code should be displayed on the screen.

- Specify the Javascript, which is code that is written to execute upon certain actions that the user of the browser may take. Javascript is similar (yet, very different) than R. This language is intended specifically to work within a browser to manipulate and change either the contend (HTML) or the look and feel (CSS) of the website upon user action.

Technically, these is a difference between a static and a dynamic website. A static website is one that does not store any information about a user's interaction with the elements on the website. When a user "visits" a static website, all they are doing is downloading the HTML, CSS, and Javascript files that will load within the user's web browser. A dynamic website, however, is much more complex. While a dynamic website must also have HTML, CSS, and Javascript, it also has *cookies* and *sessions*. A *session* is a collection of information about the user that the web server

keeps track of. Instead of the browser just simply downloading a website's code, the user can manipulate things within the website which would also be reflected on a web server. Examples of this include the more recent cloud-based software, as well as any website that allows you to "login". Such websites are built in a very complex coding environment, and often use either Perl, Python, Java, C#, Ruby, PHP, and more recently (although not very popular) R. We do not need to dive into these details, and we will work primarily on the *client-side* of web sites.

To understand the structure of a website, we will explore an example that is attached to your module. The example is a very simple example, simple enough for us to dissect line by line:

```
1  <!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.1//EN" "http://www.w3.org/TR/xhtml11/
       DTD/xhtml11.dtd">
2  <html xmlns="http://www.w3.org/1999/xhtml">
3  <head>
4  <meta http-equiv="Content-Type" content="text/html; charset=utf-8" />
5  <link rel="stylesheet" type="text/css" href="style.css" />
6
7  <title>Reverted by Bryant Smith</title>
8  </head>
9
10 <body>
11     <div id="page">
12
13        <div id="header">
14          <h1>Reverted    </h1>
15            <h2>The first XHTML 1.1 Validated Template by Bryant Smith.</h2>
16
17       </div>
18         <div id="bar">
19            <ul>
20                <li><a href="#">Home</a></li>
21                <li><a href="#">Pictures</a></li>
22                <li><a href="#">My Bio</a></li>
23                <li><a href="#">My Company</a></li>
24                <li><a href="#">Portfolio</a></li>
25                <li><a href="#">Contact</a></li>
26            </ul>
27       </div>
28         <div class="contentTitle"><h1>Template Usage</h1></div>
29         <div class="contentText">
30          <p>You may use this template on any site, anywhere, for free just<strong
   > please leave the link back to me in the footer</strong>. This template
   validates XHTML 1.1, CSS Validates as well; enjoy :) </p>
31           <p>Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Integer mi.
    Vivamus sit amet neque vitae sapien bibendum sodales. Curabitur elementum.
   Duis imperdiet. Donec eleifend porttitor sapien. Praesent leo. Quisque auctor
   velit sed tellus. Suspendisse potenti. Aenean laoreet imperdiet nunc. Donec
   commodo suscipit dolor. Aenean nibh. Sed id odio. Aliquam lobortis risus ut
   felis. Sed vehicula pellentesque quam.</p>
32           <p>Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Integer mi.
    Vivamus sit amet neque vitae sapien bibendum sodales. Curabitur elementum.
   Duis imperdiet. Donec eleifend porttitor sapien. Praesent leo. Quisque auctor
   velit sed tellus. Suspendisse potenti. Aenean laoreet imperdiet nunc. Donec
   commodo suscipit dolor. Aenean nibh. Sed id odio. Aliquam lobortis risus ut
   felis. Sed vehicula pellentesque quam.</p>
```

```
33          <p><a href="index.html">     (read more)</a></p>
34        </div>
35        <div class="contentTitle">
36          <h1>Gettysburg Address</h1>
37        </div>
38        <div class="contentText">
39          <p>Four score and seven years ago our fathers brought forth on this
     continent, a new nation, conceived in Liberty, and dedicated to the proposition
      that all men are created equal.</p>
40
41 <p>Now we are engaged in a great civil war, testing whether that nation, or any
      nation so conceived and so dedicated, can long endure. We are met on a great
      battle-field of that war. We have come to dedicate a portion of that field, as
      a final resting place for those who here gave their lives that that nation
      might live. It is altogether fitting and proper that we should do this.</p>
42
43 <p>But, in a larger sense, we can not dedicate -- we can not consecrate -- we can
      not hallow -- this ground. The brave men, living and dead, who struggled here,
      have consecrated it, far above our poor power to add or detract. The world will
       little note, nor long remember what we say here, but it can never forget what
      they did here. It is for us the living, rather, to be dedicated here to the
      unfinished work which they who fought here have thus far so nobly advanced. It
      is rather for us to be here dedicated to the great task remaining before us --
      that from these honored dead we take increased devotion to that cause for which
       they gave the last full measure of devotion -- that we here highly resolve
      that these dead shall not have died in vain -- that this nation, under God,
      shall have a new birth of freedom -- and that government of the people, by the
      people, for the people, shall not perish from the earth.</p>
44          <p><a href="index.html"> (read more)</a></p>
45        </div>
46        <div class="contentTitle"><h1>Yet Another?</h1></div>
47        <div class="contentText">Each title is an H1 tag, so choose them carefully
     :)</div>
48 </div>
49        <div id="footer"><a href="http://www.bryantsmith.com">web page designer </
     a> <a href="http://www.bryantsmith.com">bryant smith</a></div>
50 </body>
51 </html>
```
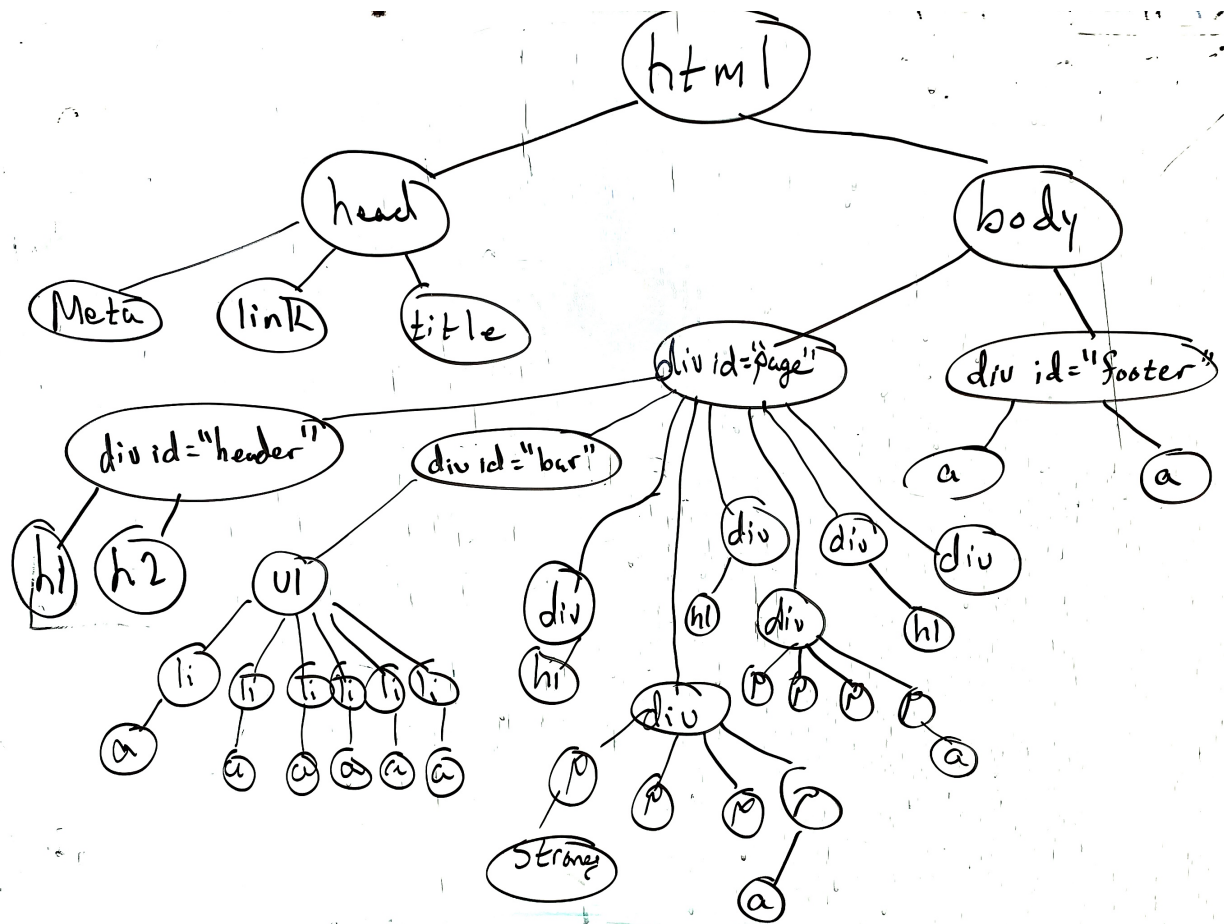
Listing 1: HTML Code for a basic static website.

Looking at the code, we notice that this looks very similar to XML. In fact, the coding structure works the same way. Information that we want to show on an HTML Web Page is determined by a tagging system, where information is placed in between tags. In the tags themselves, attributes can be specified with an attribute name = value pair. The tags that HTML can use are a fixed collection (the entire collection and meaning of them can be found here: https://www.w3schools.com/html/. Some of the the tags are themselves embedded in other tags (just as we had seen with general XML). Here are the first few tag meanings:

- <html></html> - all content on the website will be embedded in between these two tags. It has one attribute, namely the xmlns attribute. This attribute tells the web browser where it can find the version of html that the page is using.

- <head></head> This is called header information, and contains mainly meta information. Most of the information in this section does not appear directly on the website.

- `<meta/>` - specifies meta information about the website. Search engines often use these tags, and you can specify certain information to help optimize the search engines (such as Google, Bing, etc). Notice this is one tag instead of two. `<meta/>` is short for writing `<meta></meta>`, and if often used when there is nothing to put in between the tags.

- `<link />` - define a relationship between this HTML document and an outside document for which information should be inserted. Usually used to load in information about a .css file or a javascript file.

- `<title></title>` - the title of the website. This is the text that appears in the tab of a browser.

- `<body></body>` - the body of the HTML document. This is where you place the information you want to display on a page.

- `<div></div>` - a section in the HTML Page. This is the primary mechanism used to organize information on a website.

- `<h1></h1>` and `<h2></h2>` - header text which has different sizes. H1 is bigger text than H2.

- `<p>` - defines a paragraph in the HTMl Page. Usually, this tag is embedded within a `<div>` tag.

- `<a></a>` - a hyperlink. One defines where the link should bring the user by setting the `href=` attribute equal to the website in quotations.

- `<ul></ul>` - an unordered list. Each element in the list is defined by the `<li></li>` tag.

Further notice in the code that many of the `<div>` tags have an attribute of either `id=` or `class=`. This is a way to uniquely define the section on the page so that one can refer to that section later. The `class=` is called the class of the section, and it specifies a type of style that the section is to follow. General rule of thumb is that there can be multiple objects that reference the same class, but only one single object on the page that has an id. Put differently, each `id=` should be unique on the page. Furthermore, the notion of embedded tags means that the HTML document has a *tree structure* of tags. In what is called the *document object model (DOM)*, the DOM treats an HTML document as a tree of element tags. Each node in the tree will have properties, specified in the attribute of the tag. The node will also have "element text", which is any text in between the tags themselves. We can see from the below figure the tree structure for the HTML above (with my extremely poor handwriting):

But how does one define *how* the elements specified in the tree appear on the screen? This is where a second language comes into play, namely the cascading style sheet language. Cascading style sheets work very simplistically. We define a *style* by specifying certain values assigned to pre-defined attributes. We can define a style based on all tags of a certain type. We can also define styles for specific ids in the HTML page. Last, styles can be defined and grouped by what are called classes. When an element in an HTML page refers to a "class name", it is referred to the type of style that that content should display. Here is the example style sheet for our example website:

```
1  /* A Free Design by Bryant Smith (bryantsmith.com) */
2
3  html, body {
4  text-align: center;
5  }
6  p {text-align: left;}
7
8  body {
9     margin: 0;
10    padding: 0;
11    background: #697281;
12    font-family:  Arial, Helvetica, sans-serif;
13    font-size: 13px;
14    color: #F1F5F8;
```

```
15    background: #26476E url(background.png) repeat-x;
16 }
17 *
18 {
19    margin: 0 auto 0 auto;
20  text-align:left;}
21
22 #page
23 {
24    margin: 0 auto 0 auto;
25    display: block;
26    height:auto;
27    position: relative;
28    overflow: hidden;
29    width: 753px;
30 }
31
32 #header
33 {
34 margin-top:25px;
35 padding-top:63px;
36 padding-left:30px;
37 padding-right:30px;
38 background-image:url(header.png);
39 background-repeat:no-repeat;
40 height:122px;
41 width:753px;
42 overflow:hidden;
43 }
44
45 #header h1, #header h2
46 {
47 margin-left:0px;
48 width:693px;
49 position:relative;
50 color:#000000;
51 text-align:center;
52 font-size:15px;
53 font-weight:bold;
54 font-family:Arial, Helvetica, sans-serif;
55 clear:both;
56 }
57
58 #header h1
59 {
60 margin-left:0px;
61 width:693px;
62 position:relative;
63 color:#26476E;
64 font-size:30px;
65 }
66
67
68
69 #bar
70 {
```

```
71  background-image:url(bar.png);
72  background-repeat:no-repeat;
73  height:53px;
74  width:750px;
75  padding-bottom:15px;
76  text-align:center;
77  margin:0 auto 0 auto;
78  }
79
80  #bar ul {
81    margin: 0px auto 0px auto;
82    padding: 0px;
83    list-style-type: none;
84    height: 74px;
85    text-align:center;
86  }
87  #bar ul li , #bar ul li a, #bar ul li a:visited{
88    margin: 0px;
89    padding: 0px;
90    height: 30px;
91    overflow: hidden;
92    display: inline;
93    text-align:center;
94    line-height:65px;
95    color:#26476E;
96    padding-left:6px;
97    padding-right:6px;
98  }
99
100 #bar ul li a:hover
101 {
102 color:#000000;
103 text-decoration:underline;
104
105 }
106
107 .contentTitle
108 {
109 width:700px;
110 height:44px;
111 margin-top:20px;
112 margin-left:0px;
113 margin-right:25px;
114 background-image:url(tab_back.png);
115 background-repeat:no-repeat;
116
117 }
118
119 .contentTitle h1
120 {
121 margin-left:21px;
122 padding-top:8px;
123 font-size:21px;
124 font-weight:bold;
125 letter-spacing:-1px;
126 color:#FFFFFF;
```

```
127 }
128
129
130 . contentText
131 {
132 width :680 px ;
133 padding - left :30 px ;
134 padding - right :30 px ;
135 font - size :13 px ;
136 color :# FFFFFF ;
137 line - height :30 px ;
138 }
139
140 a
141 {
142 text - decoration : none ;
143 color :# F1F3F8 ;
144 }
145
146 a : hover
147 {
148 border - bottom :1 px dotted # F1F3F8 ;
149 }
150
151
152 # footer {
153 width :  750 px ;
154 height :20 px ;
155 text - align : center ;
156 font - size :9 px ;
157 color :# 223 F8D ;
158 padding - top :20 px ;
159 }
160
161 html , body {
162 text - align :  center ;
163 }
164 p {
165 margin :7 px ;
166 text - align :  left ;
167 text - align : justify ;
168 }
```

Listing 2: The CSS document for the example website.

The style sheet tells the browser which will read the HTML document how to display each component in the tree. For example, the body tag will have the following attributes:

```
1 body {
2   margin :  0;
3   padding :  0;
4   background :  #697281;
5   font - family :   Arial ,  Helvetica ,  sans - serif ;
6   font - size :  13 px ;
7   color :  # F1F5F8 ;
8   background :  #26476 E url ( background . png ) repeat - x ;
```

```
9 }
```

Listing 3: How a browser should display the body tag.

Notice that each attribute and value is assigned as follows: `attribute name : attribute value`. Each pair is followed by a semi-colon, and all of the pairs are in a block of braces. The attributes will apply to all elements in the body, unless another style that applies to only a sub-element has been defined. A complete guide to learning CSS can be found at https://www.w3schools.com/css/. Generally, styles are defined for one of four types before the braces:

- \* - means apply the style to all elements on the page.

- .classname - means to define a class with the name "classname"

- tag - means to define the style for all tags with the name "tag"

- #id - means to apply the style to any element with the id #id.

For example, the `<div></div>` block, with the id equal to "bar", in the HTML document will have the following style:

```
1  #bar
2  {
3  background-image:url(bar.png);
4  background-repeat:no-repeat;
5  height:53px;
6  width:750px;
7  padding-bottom:15px;
8  text-align:center;
9  margin:0 auto 0 auto;
10 }
```

Listing 4: The style of the div element which has the id name "bar".

Likewise, all tags the reference the class "contentText", like some of our `<div>` blocks, will have the following style:

```
1  .contentText
2  {
3  width:680px;
4  padding-left:30px;
5  padding-right:30px;
6  font-size:13px;
7  color:#FFFFFF;
8  line-height:30px;
9  }
```

Listing 5: The style for a class with a given name we defined as contentText.

At times, we want to change the style of a tag, but only if that tag is a child element in the DOM tree of another tag. For example, we may want to change the style of an unordered list (tag `<ul>`), but only those that are children elements of the element with id "bar". Hence, we would have:

```
1  #bar ul {
2    margin: 0px auto 0px auto;
3    padding: 0px;
4    list-style-type: none;
5    height: 74px;
6    text-align:center;
7  }
```

Listing 6: The style for only the ul element child of the div element with id bar.

Note that if we were to add another unordered list which is a child of a different element on the page, then this style would **not** apply to that ul, since it would not be a child of an element with the id name "bar". Suffice to say, all information on a website can hence be ascertained from its *page source*, or in other words, its HTML. In order to retrieve information from an HTML page, we can navigate the HTML tree (also called the DOM Tree). In addition, we have one additional component here, which is how to handle user interaction on a website. What happens if a user clicks, scrolls, types, etc on the web page? Ideally, we want to "listen" for these types of actions, call a function in response to it, and potentially change something in the DOM tree if the user did something on an element we expected themt to do something (that is, change the HTML code on the fly). This is where Javascript comes into play.

Javascript (which has absolutely no relation to the language Java) is a programming language that is intended to change information (HTML) and how it is displayed (CSS) in a browser by listening to a user's actions (clicking, scrolling, typing, etc). Essentially, we define a Javascript function, and code it to take an input (which is usually the action of a user) and change something in the DOM Tree (or the style). Each element in the tree is then assigned a "listener". If a listener is assigned to an element, then depending on what the user does to that element (hovers over it, scrolls over it, clicks on it, clicks off it, etc), a corresponding pre-defined function will be called.

Let us take an example. Suppose we wanted to change the text inside the h1 tags inside the `div` tag with the id "header" to "HOVERED!", any time the user hovers the mouse pointer over the `div` box itself. We can add this functionality by creating two functions in Javascript added in the header of the HTML file:

```
1  <head>
2  <meta http-equiv="Content-Type" content="text/html; charset=utf-8" />
3  <link rel="stylesheet" type="text/css" href="style.css" />
4  <script>
5    function changeTextOn(){
6      var the_header = document.getElementById("header");
7      var the_hl_tag = the_header.childNodes[1];
8      the_hl_tag.innerHTML = "HOVERED ";
9    }
10   function changeTextOff(){
11     var the_header = document.getElementById("header");
12     var the_hl_tag = the_header.childNodes[1];
13     the_hl_tag.innerHTML = "Reverted ";
14   }
15 </script>
16 <title>Reverted by Bryant Smith</title>
17 </head>
```

Listing 7: Changing the head tag content to include a script tag which is where the Javascript code will go to add the ability to change information on the page based on what the user does

on the website.

We first define the `changeTextOn()` function, which will first find the element in the tree with the id "header". Once it does that, it will look at all the child nodes. Lists (also called arrays in Javascript) are indexed starting at 0. The child element at index 1 is the `h1` tag. We then ask it to take the inner HTML in this tag and assign it to the word "HOVERED". The other function does the same, except with the word "Reverted". Notice how easy it is to manipulate the data on the HTML by using various functions and variables in javascript, thinking about the HTML page as a "tree" of tags. The functions will be called based on *events*. We can tell the browser to "listen" to the div box for the mouse to hover over it, as well as leave it, by using the `onmouseover` property, and assigning it to the function we just defined, as well as using the `onmouseout` property for when the user hovers off the `div` box, and assigning this to the other function name:

```
1 <div id="header" onmouseover="changeTextOn()" onmouseout="changeTextOff()">
2          <h1>Reverted </h1>
3            <h2>The first XHTML 1.1 Validated Template by Bryant Smith.</h2>
4
5 </div>
```

Listing 8: Adding event listeners to the element div.

Make these changes to the file `index.html` by opening it in a text editor (Notepad in Windows, Textedit in Mac), save the html file, and load it in your browser by double clicking on it. See what happens when you hover over the box! For more information on basic Javascript, visit https://www.w3schools.com/js/.

**Exercises**

7. Use the Javascript documentation referenced above to write a function in Javascript that will return all hyperlinks that are on the test webpage `index.html`.

## 4  Tutorial 2: Navigating an HTML Page via R

Now that we have a basic understanding of HTML, CSS, JS, and DOM, as well as some of the functions in JS that are used to find and manipulate elements in HTML, we can leverage the packages in R to do just this with our HTML Document. When we scrape the web, we are essentially parsing through an HTML Document to determine which information is relevant and which information is not. In order to scrape information, we typically need to understand the structure of the HTML website, depending on our scraping task. Hence, in this next portion of the tutorial, we will leverage R to navigate through an HTML document.

When we are performing web scraping on static websites, luckily for us, R has an easy to use package that only really necessitates us in having a mild understanding of the HTML DOM Tree for the site. The package we will use here is `rvest`. This package is a wrapper (that is, it calls other functions) around more complicated packages that allow us to traverse the HTML. Let us parse our own website as an example. Move the `index.html` file to your R working directory (or, set your R working directory to the directory in which `index.html` is located). Next, install the `rvest` package and load it. We will load the html document into R as such:

```
1 > library(rvest)
2 >
```

```
3  > the_html<-read_html("index.html")
4  > the_html
5  {xml_document}
6  <html xmlns="http://www.w3.org/1999/xhtml">
7  [1] <head>\n<meta http-equiv= ...
8  [2] <body>\r\n    <div id="pa ...
```

Listing 9: Loading in the HTML file in R for processing.

Now the object that is stored in the variable `the_html` has a lot of structure that we can leverage the extract any information we may need. For example, we can start by understanding its basic tree structure:

```
1   > html_structure(the_html,indent=5)
2   <html [xmlns]>
3        <head>
4             <meta [http-equiv, content]>
5             <link [rel, type, href]>
6             <script>
7                  {cdata}
8             <title>
9                  {text}
10       <body>
11            {text}
12            <div#page>
13                 {text}
14                 <div#header [onmouseover, onmouseout]>
15                      {text}
16                      <h1>
17                           {text}
18                      {text}
19                      <h2>
20                           {text}
21                      {text}
22                 {text}
23                 <div#bar>
24                      {text}
25                      <ul>
26                           <li>
27                                <a [href]>
28                                     {text}
29                           {text}
30                           <li>
31                                <a [href]>
32                                     {text}
33                           {text}
34                           <li>
35                                <a [href]>
36                                     {text}
37                           {text}
38                           <li>
39                                <a [href]>
40                                     {text}
41                           {text}
42                           <li>
43                                <a [href]>
```

```
44                                      {text}
45                          {text}
46                          <li>
47                                  <a [href]>
48                                          {text}
49                          {text}
50                  {text}
51                  <div.contentTitle>
52                      <h1>
53                          {text}
54                  {text}
55                  <div.contentText>
56                      {text}
57                      <p>
58                          {text}
59                      <strong>
60                              {text}
61                      {text}
62                  {text}
63                  <p>
64                          {text}
65                  {text}
66                  <p>
67                          {text}
68                  {text}
69                  <p>
70                      <a [href]>
71                              {text}
72                  {text}
73              {text}
74              <div.contentTitle>
75                  {text}
76                  <h1>
77                          {text}
78                  {text}
79              {text}
80              <div.contentText>
81                  {text}
82                  <p>
83                          {text}
84                  {text}
85                  <p>
86                          {text}
87                  {text}
88                  <p>
89                          {text}
90                  {text}
91                  <p>
92                      <a [href]>
93                              {text}
94                  {text}
95              {text}
96              <div.contentTitle>
97                  <h1>
98                          {text}
99              {text}
```

```
100                  <div.contentText >
101                       {text}
102                  {text}
103            {text}
104            <div#footer >
105                  <a [href]>
106                       {text}
107                  {text}
108                  <a [href]>
109                       {text}
110            {text}
```

Listing 10: Getting the tree structure of an HTML document.

Now suppose that we would like to get all of the text that is stored in all of the paragraph tags. We can very easily do this by first finding all elements that are paragraph tags, and then subsequently by converting this list to text:

```
1 > p_tags<-html_nodes(the_html,"p")
2 > p_tags
3 {xml_nodeset (8)}
4 [1] <p>You may use this template on any site, anywhere, for fr ...
5 [2] <p>Lorem ipsum dolor sit amet, consectetuer adipiscing eli ...
6 [3] <p>Lorem ipsum dolor sit amet, consectetuer adipiscing eli ...
7 [4] <p><a href="index.html">  (read more)</a></p>
8 [5] <p>Four score and seven years ago our fathers brought fort ...
9 [6] <p>Now we are engaged in a great civil war, testing whethe ...
10 [7] <p>But, in a larger sense, we can not dedicate -- we can n ...
11 [8] <p><a href="index.html"> (read more)</a></p>
12 > p_text<-html_text(p_tags)
13 > p_text
14 [1] "You may use this template on any site, anywhere, for free just please leave
       the link back to me in the footer. This template validates XHTML 1.1, CSS
       Validates as well; enjoy :) "
15 [2] "Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Integer mi. Vivamus
        sit amet neque vitae sapien bibendum sodales. Curabitur elementum. Duis
       imperdiet. Donec eleifend porttitor sapien. Praesent leo. Quisque auctor velit
       sed tellus. Suspendisse potenti. Aenean laoreet imperdiet nunc. Donec commodo
       suscipit dolor. Aenean nibh. Sed id odio. Aliquam lobortis risus ut felis. Sed
       vehicula pellentesque quam."
16 [3] "Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Integer mi. Vivamus
        sit amet neque vitae sapien bibendum sodales. Curabitur elementum. Duis
       imperdiet. Donec eleifend porttitor sapien. Praesent leo. Quisque auctor velit
       sed tellus. Suspendisse potenti. Aenean laoreet imperdiet nunc. Donec commodo
       suscipit dolor. Aenean nibh. Sed id odio. Aliquam lobortis risus ut felis. Sed
       vehicula pellentesque quam."
17 [4] "  (read more)"
18 [5] "Four score and seven years ago our fathers brought forth on this continent, a
        new nation, conceived in Liberty, and dedicated to the proposition that all
       men are created equal."
19 [6] "Now we are engaged in a great civil war, testing whether that nation, or any
       nation so conceived and so dedicated, can long endure. We are met on a great
       battle-field of that war. We have come to dedicate a portion of that field, as
       a final resting place for those who here gave their lives that that nation
       might live. It is altogether fitting and proper that we should do this."
```

```
20 [7] "But , in a larger sense , we can not dedicate -- we can not consecrate -- we
        can not hallow -- this ground. The brave men , living and dead , who struggled
        here , have consecrated it , far above our poor power to add or detract. The
        world will little note , nor long remember what we say here , but it can never
        forget what they did here. It is for us the living , rather , to be dedicated
        here to the unfinished work which they who fought here have thus far so nobly
        advanced. It is rather for us to be here dedicated to the great task remaining
        before us -- that from these honored dead we take increased devotion to that
        cause for which they gave the last full measure of devotion -- that we here
        highly resolve that these dead shall not have died in vain -- that this nation ,
         under God , shall have a new birth of freedom -- and that government of the
        people , by the people , for the people , shall not perish from the earth."
21 [8] " (read more)"
```

Listing 11: Extracting text from all the paragraph tags.

Now suppose we wanted to extract the URLs to all hyperlinks on the webpage, then we can do this by using:

```
1 > a_tags<-html_nodes(the_html ,"a")
2 > a_tags
3 {xml_nodeset (10)}
4  [1] <a href ="#">Home </a>
5  [2] <a href ="#">Pictures </a>
6  [3] <a href ="#">My Bio </a>
7  [4] <a href ="#">My Company </a>
8  [5] <a href ="#">Portfolio </a>
9  [6] <a href ="#">Contact </a>
10  [7] <a href ="index.html">  (read more)</a>
11  [8] <a href ="index.html"> (read more)</a>
12  [9] <a href ="http ://www.bryantsmith.com">web page designer </a>
13 [10] <a href ="http ://www.bryantsmith.com">bryant smith </a>
14 > html_attr(a_tags ,"href")
15  [1] "#"                      "#"
16  [3] "#"                      "#"
17  [5] "#"                      "#"
18  [7] "index.html"             "index.html"
19  [9] "http ://www.bryantsmith.com" "http ://www.bryantsmith.com"
```

Listing 12: Extract the attribute "href" from all "a" tags.

**Exercises**

8. Loop up the documentation for the rvest package. Write code in R to extract only the text between <a> tags that are only in the div tag with id page. Ensure the code will not explore tags that are not under this tag.

# 5  Tutorial 3: Static-Based Web Scraping

The fundamental difference between static and dynamic web scraping lies in the design of the website itself. In some websites (such as Wix-Generated Websites, for example), the "website" is not actually the HTML that you download. Hence, it would be impossible to obtain data that is not in an HTML document yet. Data that is only generated due to calling a Javascript

function, one that cannot be invoked without a browser due to the complexity. However, in other situations, the website (despite being a dynamic website) may be simple enough where you will not need to resort to more advanced tools.

Put differently, if you are able to download the HTML document, and your target data is in the HTML document, you can take a *static* approach to web scraping. This approach is straightforward: download the HTML page, parse it with a tool (such as `rvest`), and save the data to your intended location. Sounds easy, right? Well, there are a few things to consider. First, you must ensure that the data itself is indeed in the HTML. If it is not, then there is no point in parsing the HTML code, since your data wont be there! You would need to resort to a dynamic approach (see below). On the other hand, if your data is in the HTML page that you downloaded, then congrats! All you need to do is parse it. However, even parsing necessitates having a background knowledge on the structure of the HTML document.

Understanding the "structure" of a website boils down to trying to understand common patterns where your intended data most likely lie. This necessitates you to essentially reverse engineer a website's HTML, if at all possible. The good news is that most websites store their information in well-constructed and predictable ways. The job of the author of the web scraper is to first understand what that structure is and to subsequently write code to exploit that structure. Let us take an example. Suppose we wanted to download county level election data in New Jersey but only during presidential elections. We happen to find that Wikipedia has a nice collection of this information for the most recent 2016 election:



This is a very nice looking table, and it would be nice if all elections had this table as well:

However, we don't want to go election by election by hand and copy and paste this table. It would be great if we can write a scrapper to download this information. We first want to understand its structure. In Google Chrome (which I recommend as your only browser for doing this), if you right click on the table (or anywhere), you will see on the bottom an option that says "inspect":



Click this, and you will have the HTML code (among other things) show up the right hand side of the screen. This tool is a very useful tool to understand how the HTML/CSS/JS is generating the content. When you hover over an element in the HTML code, it will highlight on your right hand side screen:

Ideally, we would like to see if there is any common pattern that can uniquely lead us to this specific table. We notice that the table does not have a unique id, but it does follow the class `wikitable sortable jquery-tablesorter`. Upon searching in the rest of the document for "wikitable", we notice a few more tables show up. However, a search for "sortable" shows us that this is the only table that uses this class. This may be beneficial for us, and we will keep this in mind:



Now we would like to see if other pages follow a similar structure. If it does, then we are in great luck, and we can leverage that. After another Google Search for the election before this one, we notice that the url for the pages itself follows an interesting pattern. Namely that the url is exactly the same for each type of page, exact with the election year replaced. Further confirming, we can see this is indeed a way to reference the different elections:

Looking at the structure again by clicking "inspect", we see that once again this table follows the same exact classname as the 2016 election table. Furthermore, a quick search of "sortable" again uniquely identifies this table. It looks like we hence have a way to determine where in the HTML tree this table lies. Do we need to confirm this for all pages? No, since if we write a quick script to see if there is only one unique table per election page, then we can take a gamble and run the scrapper, then later clean the data by hand:



We also by going back through the election years that the first year which shows county-level data is the 1924 election:

Hence, we have a clear objective in mind for the data we would like to scrape: All County-Level Election Results for only Democrats and Republicans in every presidential election year starting from the 1924 election and ending at the 2016 election. Now it should be noted that upon further investigation, some tables have the class name `wikitable sortable jquery-tablesorter` while others just have `wikitable sortable`. Either way, given that we have already established that it is likely all tables uniquely have the class name "sortable" in it, we can hence refer to this table using a CSS Selectors (see https://www.w3schools.com/cssref/css$_s$electors.asp for more on this). The selector is a way to uniquely "select" the element. We start with the tag name, followed by the class name, and all subsequent classes, with periods in between each. So if we want to select a table that is in the class wikitable and sortable, we would indicate this by "table.wikitable.sortable".

With all of this in mind, let us test how to extract this from the 2016 election page. First, we will insert the url into a variable and pass it into the read html function we discussed earlier. The package rvest will parse the HTML into a tree, which makes it easy for us to navigate. Once we do this we can now extract the table by id. Last, once we are brought to the table element in the HTML tree, we need to extract the data. Traditionally, this has been a pain. Luckily, rvest has a nice easy to use function to accomplish this, namely the `html_table` function. This will extract the data from the HTML table and store it in a data frame. For now, we will simply just store the results as a csv file for later processing:

```
1 > the_html<-read_html("https://en.wikipedia.org/wiki/2016_United_States_
     presidential_election_in_New_Jersey")
2 >
3 > county_html_table<-html_nodes(the_html,"table.wikitable.sortable")
4 >
5 > county_data<-html_table(county_html_table)
6 >
7 > county_data
8 [[1]]
9         County Clinton votes Clinton % Trump votes Trump % Other votes  Other %
```

```
10 1     Atlantic      60,924   51.0%    52,690   44.1%    3,677   3.13%
11 2       Bergen     231,211   54.2%   175,529   41.1%   12,556   2.99%
12 3   Burlington     121,725   54.2%    89,272   39.7%    7,946   3.63%
13 4       Camden     146,717   63.4%    72,631   31.4%    7,244   3.20%
14 5     Cape May      18,750   37.5%    28,446   57.0%    1,526   3.13%
15 6   Cumberland      27,771   50.4%    24,453   44.4%    1,780   3.30%
16 7        Essex     240,837   76.2%    63,176   20.0%    6,921   2.23%
17 8    Gloucester     66,870   46.9%    67,544   47.4%    5,128   3.67%
18 9       Hudson     163,917   73.2%    49,043   21.9%    6,415   2.92%
19 10    Hunterdon     28,898   39.7%    38,712   53.2%    3,226   4.55%
20 11       Mercer    104,775   65.6%    46,193   28.9%    5,561   3.55%
21 12    Middlesex    193,044   58.0%   122,953   37.0%   10,105   3.10%
22 13     Monmouth    137,181   42.3%   166,723   51.5%   10,473   3.33%
23 14       Morris    115,249   44.9%   126,071   49.1%    9,096   3.63%
24 15        Ocean     87,150   31.1%   179,079   63.9%    8,133   2.96%
25 16      Passaic    116,759   58.8%    72,902   36.7%    5,141   2.64%
26 17        Salem     11,904   39.6%    16,381   54.4%    1,209   4.10%
27 18     Somerset     85,689   53.4%    65,505   40.8%    5,898   3.75%
28 19       Sussex     24,212   32.0%    46,658   61.8%    3,256   4.39%
29 20        Union    147,414   66.4%    68,114   30.7%    6,447   2.90%
30 21       Warren     17,281   34.3%    29,858   59.2%    2,097   4.26%
```

Listing 13: Extracting the Election County Table Results.

We now can see that our code works for a single year, but what about other years? Let us try it out on all election years from 1924 to 2016. We can see that the url for the specific election year is exactly the same with the exception of the year. Hence, we can generate a vector of years, starting at 1924, ending at 2016, and going up by every 4 years. Then, we can write a loop to iterate for each year to insert the year in the URL, run the same code we just ran above but for that year, and save the result in a csv table (although, we will save it in a list for now). Let us see how this goes:

```
1  > election_years<-seq(1924,2016,4)
2  > county_data_year<-list()
3  > for(year in election_years){
4  +    the_url<-paste("https://en.wikipedia.org/wiki/",year,"_United_States_
        presidential_election_in_New_Jersey",sep="")
5  +
6  +    the_html<-read_html(the_url)
7  +
8  +    county_html_table<-html_nodes(the_html,"table.wikitable.sortable")
9  +
10 +    county_data<-html_table(county_html_table)
11 +
12 +    county_data_year[[paste(year,sep="")]] = county_data
13 + }
14 Error: Table has inconsistent number of columns. Do you want fill = TRUE?
```

Listing 14: Extracting the results for multiple years but receiving an error in doing so.

Ah, so, we now encounter an error. What are we to do? Sometimes when we scrape, errors are okay, and we just want to ignore them, since they act as just annoying signals to us that we do not care about. Other times, errors may indicate that we wrote something wrong in our code, and we hence need to *debug* the code. In our case, we are not sure, so let us just debug to see

what the deal is. One method for debugging in loops is to see how far the loop got before it broke down. In the code below, we can see that the code worked just fine up until the year 1964, which is where it broke down. Hence, we will assign the variable year to the value 1964 and see why it broke down by running each line inside the loop. Doing so leads us to the error when we try to convert the element to the data frame. That is, we can see the html is downloaded just fine, that the HTML Tree is generated just fine, that the table is found just fine, but the error lies in the conversion when we call the function `html_table`:

```
1 > names(county_data_year)
2  [1] "1924" "1928" "1932" "1936" "1940" "1944" "1948" "1952" "1956" "1960"
3 > year <- 1964
4 > the_url<-paste("https://en.wikipedia.org/wiki/",year,"_United_States_
    presidential_election_in_New_Jersey",sep="")
5 >
6 >   the_html<-read_html(the_url)
7 >
8 >   county_html_table<-html_nodes(the_html,"table.wikitable.sortable")
9 > county_html_table
10 {xml_nodeset (1)}
11 [1] <table width="70%" class="wikitable sortab ...
12 > county_data<-html_table(county_html_table)
13 Error: Table has inconsistent number of columns. Do you want fill = TRUE?
14 > ?html_table
```

Listing 15: Debugging the problem and trying to find where the error started.

Upon looking at the documentation, we see that there is an input into this function indicated in the documentation: "fill If TRUE, automatically fill rows with fewer than the maximum number of columns with NAs." So our solution may be as simple as just setting this value equal to true in the input. We will do this. Furthermore, as a sanity check, let us see if there is anything is being stored in the list in each iteration. We do this by printing the element:

```
1 > election_years<-seq(1924,2016,4)
2 > county_data_year<-list()
3 > for(year in election_years){
4 +   the_url<-paste("https://en.wikipedia.org/wiki/",year,"_United_States_
    presidential_election_in_New_Jersey",sep="")
5 +
6 +   the_html<-read_html(the_url)
7 +
8 +   county_html_table<-html_nodes(the_html,"table.wikitable.sortable")
9 +   print(year)
10 +   print(county_html_table)
11 +   county_data<-html_table(county_html_table,fill=TRUE)
12 +
13 +   county_data_year[[paste(year,sep="")]] <- county_data
14 + }
15 [1] 1924
16 {xml_nodeset (1)}
17 [1] <table width="75%" class="wikitable sortable" style="text ...
18 [1] 1928
19 {xml_nodeset (0)}
20 [1] 1932
```

```
21  {xml_nodeset (0)}
22  [1] 1936
23  {xml_nodeset (0)}
24  [1] 1940
25  {xml_nodeset (0)}
26  [1] 1944
27  {xml_nodeset (0)}
28  [1] 1948
29  {xml_nodeset (0)}
30  [1] 1952
31  {xml_nodeset (0)}
32  [1] 1956
33  {xml_nodeset (0)}
34  [1] 1960
35  {xml_nodeset (1)}
36  [1] <table width="60%" class="wikitable sortable" style="text ...
37  [1] 1964
38  {xml_nodeset (1)}
39  [1] <table width="70%" class="wikitable sortable" style="text ...
40  [1] 1968
41  {xml_nodeset (1)}
42  [1] <table class="wikitable sortable" style="text-align:cente ...
43  [1] 1972
44  {xml_nodeset (1)}
45  [1] <table width="60%" class="wikitable sortable" style="text ...
46  [1] 1976
47  {xml_nodeset (1)}
48  [1] <table width="60%" class="wikitable sortable" style="text ...
49  [1] 1980
50  {xml_nodeset (1)}
51  [1] <table width="65%" class="wikitable sortable" style="text ...
52  [1] 1984
53  {xml_nodeset (1)}
54  [1] <table width="65%" class="wikitable sortable" style="text ...
55  [1] 1988
56  {xml_nodeset (1)}
57  [1] <table width="65%" class="wikitable sortable" style="text ...
58  [1] 1992
59  {xml_nodeset (1)}
60  [1] <table width="65%" class="wikitable sortable" style="text ...
61  [1] 1996
62  {xml_nodeset (1)}
63  [1] <table width="65%" class="wikitable sortable" style="text ...
64  [1] 2000
65  {xml_nodeset (1)}
66  [1] <table class="wikitable sortable" style="text-align:right ...
67  [1] 2004
68  {xml_nodeset (1)}
69  [1] <table class="wikitable sortable" style="text-align:right ...
70  [1] 2008
71  {xml_nodeset (1)}
72  [1] <table width="60%" class="wikitable sortable"><tbody>\n<t ...
73  [1] 2012
74  {xml_nodeset (1)}
75  [1] <table class="wikitable sortable"><tbody>\n<tr>\n<th>Coun ...
76  [1] 2016
```

```
77 {xml_nodeset (1)}
78 [1] <table class="wikitable sortable" style="text-align:right ...
```

Listing 16: Further debugging and partial fixing of the problem.

As we can see from the output, it appears that the 1924 data point is an outlier. All years in between 1924 and 1960 appear to not have a table of county level results. As such, we will need to sacrifice our pursuit of this data and limit our download to only what is available. Namely, all county level data from 1960 to 2016 (still not a terrible data set!). Running the code again, we get a clean data set. Next is to confirm this. Unfortunately, each entry in the list is itself a list which contains a data frame (I know, this can get confusing!). As such, we need to extract each data frame and put it purely in the list. After, we inspect each data frame from each year:

```
1  > unlist(lapply(county_data_year,is.data.frame))
2   1960   1964   1968   1972   1976   1980   1984   1988   1992   1996   2000   2004
3  FALSE  FALSE  FALSE  FALSE  FALSE  FALSE  FALSE  FALSE  FALSE  FALSE  FALSE  FALSE
4   2008   2012   2016
5  FALSE  FALSE  FALSE
6  > y<-lapply(county_data_year,function(x){x[[1]]})
7  > unlist(lapply(y,is.data.frame))
8  1960 1964 1968 1972 1976 1980 1984 1988 1992 1996 2000 2004 2008 2012
9  TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
10 2016
11 TRUE
12 >
13 > lapply(y,function(x){print(x[1:3,1:4])})
14           John Fitzgerald KennedyDemocratic
15 1    County                               #
16 2 Atlantic                           36,129
17 3   Bergen                          156,165
18   John Fitzgerald KennedyDemocratic.1 Richard Milhous NixonRepublican
19 1                                   %                               #
20 2                               46.94%                          39,158
21 3                               40.90%                         224,969
22     County Lyndon Baines JohnsonDemocratic
23 1    County                               #
24 2 Atlantic                           50,945
25 3   Bergen                          234,849
26   Lyndon Baines JohnsonDemocratic.1 Barry Morris GoldwaterRepublican
27 1                                 %                               #
28 2                             65.30%                          25,626
29 3                             59.69%                         157,899
30           Richard Milhous NixonRepublican
31 1    County                               #
32 2 Atlantic                           32,807
33 3   Bergen                          224,911
34   Richard Milhous NixonRepublican.1
35 1                                 %
36 2                             42.15%
37 3                             54.45%
38   Hubert Horatio Humphrey Jr.Democratic
39 1                                     #
40 2                                35,581
41 3                               162,182
42           Richard Milhous NixonRepublican
```

```
43  1    County                                             #
44  2 Atlantic                                           45,667
45  3   Bergen                                           285,458
46    Richard Milhous NixonRepublican.1 George Stanley McGovernDemocratic
47  1                                      %                                        #
48  2                                   59.54%                                   28,203
49  3                                   65.34%                                  147,155
50            Gerald FordRepublican Gerald FordRepublican.1
51  1    County                        #                           %
52  2 Atlantic                      36,733                       45.56%
53  3   Bergen                      237,331                      55.86%
54    Jimmy CarterDemocratic
55  1                      #
56  2                  41,965
57  3                 180,738
58            Ronald Wilson ReaganRepublican
59  1    County                                      #
60  2 Atlantic                                     37,973
61  3   Bergen                                     232,043
62    Ronald Wilson ReaganRepublican.1 James Earl CarterDemocratic
63  1                                      %                                        #
64  2                                   49.83%                                   31,286
65  3                                   55.89%                                  139,474
66            Ronald Wilson ReaganRepublican
67  1    County                                      #
68  2 Atlantic                                     49,158
69  3   Bergen                                     268,507
70    Ronald Wilson ReaganRepublican.1 Walter Fritz MondaleDemocratic
71  1                                      %                                        #
72  2                                   59.33%                                   33,240
73  3                                   63.22%                                  155,039
74            George Herbert Walker BushRepublican
75  1    County                                         #
76  2 Atlantic                                       44,748
77  3   Bergen                                       226,885
78    George Herbert Walker BushRepublican.1
79  1                                      %
80  2                                   56.33%
81  3                                   58.19%
82    Michael Stanley DukakisDemocratic
83  1                                   #
84  2                               34,047
85  3                              160,655
86            William Jefferson ClintonDemocratic
87  1    County                                         #
88  2 Atlantic                                       39,633
89  3   Bergen                                       171,104
90    William Jefferson ClintonDemocratic.1
91  1                                      %
92  2                                   43.89%
93  3                                   42.44%
94    George Herbert Walker BushRepublican
95  1                                         #
96  2                                     34,279
97  3                                    178,223
98            William Jefferson ClintonDemocratic
```

```
99  1    County                                           #
100 2  Atlantic                                      44,434
101 3    Bergen                                     191,085
102   William Jefferson ClintonDemocratic.1 Robert Joseph DoleRepublican
103 1                                      %                            #
104 2                                  53.15%                       29,538
105 3                                  52.66%                      141,164
106       County Gore votes Gore % Bush votes
107 1    Atlantic     52,880  58.0%     35,593
108 2      Bergen    202,682  55.3%    152,731
109 3 Burlington     99,506  56.1%     72,254
110       County Kerry % Kerry # Bush %
111 1    Atlantic    52.5%  55,746  46.6%
112 2      Bergen    51.9% 207,666  47.4%
113 3 Burlington    53.1% 110,411  46.1%
114          X1       X2      X3       X4
115 1    County Obama % Obama # McCain %
116 2  Atlantic   56.9%  67,830    41.8%
117 3    Bergen   54.2% 225,367    44.7%
118       County Obama%  Obama# Romney%
119 1    Atlantic 57.96%  65,600  41.10%
120 2      Bergen 55.20% 212,754  43.87%
121 3 Burlington 58.53% 126,377  40.48%
122       County Clinton votes Clinton % Trump votes
123 1    Atlantic         60,924     51.0%      52,690
124 2      Bergen        231,211     54.2%     175,529
125 3 Burlington        121,725     54.2%      89,272
126 $ '1960'
127          John Fitzgerald KennedyDemocratic
128 1    County                              #
129 2  Atlantic                         36,129
130 3    Bergen                        156,165
131   John Fitzgerald KennedyDemocratic.1 Richard Milhous NixonRepublican
132 1                                    %                             #
133 2                                46.94%                        39,158
134 3                                40.90%                       224,969
135
136 $ '1964'
137     County Lyndon Baines JohnsonDemocratic
138 1    County                              #
139 2  Atlantic                         50,945
140 3    Bergen                        234,849
141   Lyndon Baines JohnsonDemocratic.1 Barry Morris GoldwaterRepublican
142 1                                  %                            #
143 2                              65.30%                       25,626
144 3                              59.69%                      157,899
145
146 $ '1968'
147          Richard Milhous NixonRepublican
148 1    County                              #
149 2  Atlantic                         32,807
150 3    Bergen                        224,911
151   Richard Milhous NixonRepublican.1
152 1                                  %
153 2                              42.15%
154 3                              54.45%
```

27/38

```
155    Hubert Horatio Humphrey Jr.Democratic
156 1                                              #
157 2                                         35,581
158 3                                        162,182
159
160 $ '1972'
161          Richard Milhous NixonRepublican
162 1    County                                   #
163 2 Atlantic                               45,667
164 3    Bergen                             285,458
165   Richard Milhous NixonRepublican.1 George Stanley McGovernDemocratic
166 1                                    %                                        #
167 2                                59.54%                                   28,203
168 3                                65.34%                                  147,155
169
170 $ '1976'
171          Gerald FordRepublican Gerald FordRepublican.1
172 1    County                    #                          %
173 2 Atlantic               36,733                    45.56%
174 3    Bergen             237,331                    55.86%
175   Jimmy CarterDemocratic
176 1                      #
177 2                 41,965
178 3                180,738
179
180 $ '1980'
181          Ronald Wilson ReaganRepublican
182 1    County                                   #
183 2 Atlantic                               37,973
184 3    Bergen                             232,043
185   Ronald Wilson ReaganRepublican.1 James Earl CarterDemocratic
186 1                                    %                                        #
187 2                                49.83%                                   31,286
188 3                                55.89%                                  139,474
189
190 $ '1984'
191          Ronald Wilson ReaganRepublican
192 1    County                                   #
193 2 Atlantic                               49,158
194 3    Bergen                             268,507
195   Ronald Wilson ReaganRepublican.1 Walter Fritz MondaleDemocratic
196 1                                    %                                        #
197 2                                59.33%                                   33,240
198 3                                63.22%                                  155,039
199
200 $ '1988'
201          George Herbert Walker BushRepublican
202 1    County                                    #
203 2 Atlantic                               44,748
204 3    Bergen                             226,885
205   George Herbert Walker BushRepublican.1
206 1                                    %
207 2                                56.33%
208 3                                58.19%
209   Michael Stanley DukakisDemocratic
210 1                                         #
```

```
211 2                                       34,047
212 3                                      160,655
213
214 $`1992`
215          William Jefferson ClintonDemocratic
216 1    County                                   #
217 2 Atlantic                              39,633
218 3    Bergen                            171,104
219   William Jefferson ClintonDemocratic.1
220 1                                       %
221 2                                    43.89%
222 3                                    42.44%
223   George Herbert Walker BushRepublican
224 1                                       #
225 2                                    34,279
226 3                                   178,223
227
228 $`1996`
229          William Jefferson ClintonDemocratic
230 1    County                                   #
231 2 Atlantic                              44,434
232 3    Bergen                            191,085
233   William Jefferson ClintonDemocratic.1 Robert Joseph DoleRepublican
234 1                                       %                            #
235 2                                    53.15%                       29,538
236 3                                    52.66%                      141,164
237
238 $`2000`
239       County Gore votes Gore % Bush votes
240 1    Atlantic    52,880  58.0%     35,593
241 2      Bergen   202,682  55.3%    152,731
242 3 Burlington    99,506  56.1%     72,254
243
244 $`2004`
245       County Kerry % Kerry # Bush %
246 1    Atlantic   52.5%  55,746  46.6%
247 2      Bergen   51.9% 207,666  47.4%
248 3 Burlington   53.1% 110,411  46.1%
249
250 $`2008`
251        X1       X2       X3       X4
252 1    County Obama % Obama # McCain %
253 2 Atlantic   56.9%  67,830    41.8%
254 3    Bergen   54.2% 225,367    44.7%
255
256 $`2012`
257      County Obama%  Obama# Romney%
258 1    Atlantic 57.96%  65,600  41.10%
259 2      Bergen 55.20% 212,754  43.87%
260 3 Burlington 58.53% 126,377  40.48%
261
262 $`2016`
263      County Clinton votes Clinton % Trump votes
264 1    Atlantic         60,924     51.0%      52,690
265 2      Bergen        231,211     54.2%     175,529
```

```
266  3 Burlington        121,725       54.2%        89,272
```

Listing 17: Correcting the data format and inspecting the structure of each table.

Unfortunately, we can see that each table for each year follows a slightly different format. At this point, manual manipulation is most likely needed. However, it is evident that we at least have downloaded the important data, and the rest is merely small data cleaning and organizing tasks.

**Exercises**

9. In many of the presidential elections, there is another table under the "Results" section on the wikipedia page. Write R code to extract, download, and save this table as a csv file with the year as the name of the file.

# 6  Tutorial 4: Dynamic-Based Web Scraping

As we illustrated in the previous section, static web scraping is made very easy when using the `rvest` package. This is due to the fact that after years of various packages and the like, parsing HTML/CSS/JS has become well established. However, there are considerations to think about when using this method. First, the most important aspect of all of this is that when you scrape a site in a static manner, you are not interacting with that site. You are merely downloading an HTML/CSS/JS file, and working on that file directly. If a change on the website has been made, this will not be reflected in the HTML file. You can test this with the HTML code from Tutorial 2. Despite the fact that our JS code changes the HTML code in the browser, it will not do so in the file. Second, in static web scraping, we are not downloading the HTML file from a browser, but rather through an application (in our instance, through R). This means that you are in no way interacting with the website.

However, the website for which you seek to interact may be so complicated that the HTML code is changed on the fly. Meaning, the changes in the HTML code will only take place in a web browser rather than in the HTML file itself. This poses some issues when scraping. Other issues to think about when we scrape via static methods:

- Your ISP may ban you from a website. Too many calls to a website can result in your ip address being placed on a "blocked" list.

- The web server itself may block your IP address.

- The web site may dynamically update the underlying HTML code, which can only be accessed in a web browser.

- The web site may have what is called session management, which keeps track of your movements from page to page (such as "logging into" a website). You may not be able to access such information unless you are "logged in".

In these instances, our methods discussed in the previous tutorial will not hold. We will either be banned from a website, or, will be unable to gain access programatically to the data which we need. So what can we do? The answer lies in what is called *dynamic web scraping*. Earlier, we used R to essentially send a GET request to a website, which in turn downloaded HTML documents for us to process, all without a browser. In dynamic web scraping, we insert

a "middle man" into this process. Instead of analyzing HTML code from a file, we instead use a web browser to send requests to websites. The browser will in turn receive HTML code. Instead of us analyzing the code directly, we instead programatically control the browser's movements through a series of scrolls, clicks, mouse movements, typing,etc. By doing so, the HTML in the browser will change. However, since we have programmatic access to the browser, we will be able to obtain the information that is needed. This method allows us to overcome many of the problems listed earlier. The reason is simple: to the ISP and the website, they cannot tell the difference between you or a computer using the web browser.

Of course, this all comes at a cost (after all, there is no such thing as a free lunch!). We sacrifice *speed* in favor of quantity and quality of data. Given that we are doing this in a web browser, we are adding one more layer of code in between us and the web server. This means that our requests/responses will slow. Despite this, we can still get really good data that we otherwise are unable to get via static web scrapping! Let us look at a small example. The way in which we connect to a browser in R is through the Selenium software. You will need the following to run it:

- The latest version of Java installed on your computer.

- The RSelenium R pacakge installed.

Setup is not difficult, but is a bit tedious. First, go to your Google Chrome browser and look at its version number (go to Settings>About Chrome). Insert that version number in the "chromever" option in the rsDriver function. R will most likely give you an error:

```r
> library(RSelenium)
> rd<-rsDriver(browser=c("chrome"),chromever="79.0.3945.117")
checking Selenium Server versions:
BEGIN: PREDOWNLOAD
BEGIN: DOWNLOAD
BEGIN: POSTDOWNLOAD
checking chromedriver versions:
BEGIN: PREDOWNLOAD
BEGIN: DOWNLOAD
BEGIN: POSTDOWNLOAD
Error in chrome_ver(chromecheck[["platform"]], chromever) :
  version requested doesnt match versions available =
     75.0.3770.90,76.0.3809.12,76.0.3809.126,76.0.3809.25,77.0.3865.10,77.0.3865.40,78.0.3904.105
```

Listing 18: Getting RSelenium up and running but getting an error message with the wrong version number.

Find the version number that is less than you chrome version, but is the highest one that is less than your chrome. In my case, this would be "78.0.3904.11", since this is the highest version less than "79.0.3945.117". Insert this as the chromever option, and try again:

```r
> rd<-rsDriver(browser=c("chrome"),chromever="78.0.3904.11")
checking Selenium Server versions:
BEGIN: PREDOWNLOAD
BEGIN: DOWNLOAD
BEGIN: POSTDOWNLOAD
```
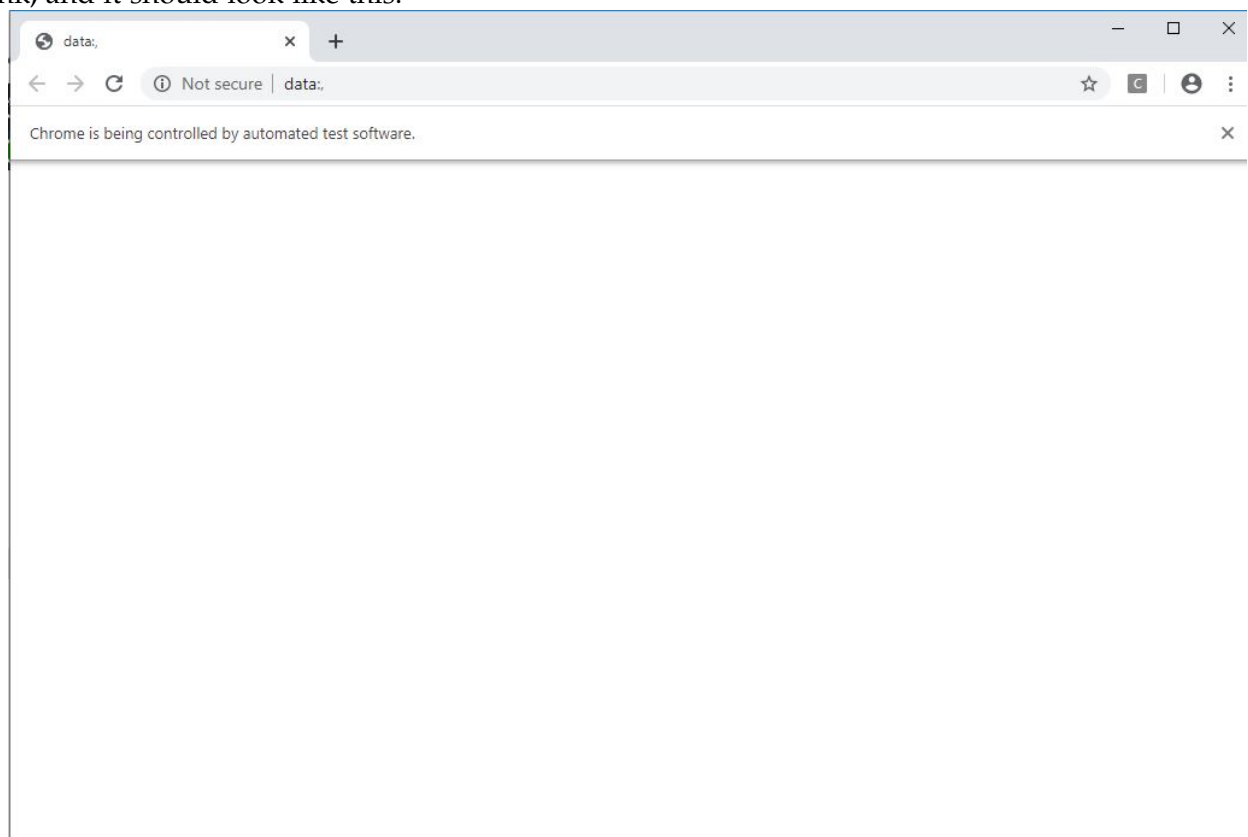
```
 6 checking chromedriver versions:
 7 BEGIN: PREDOWNLOAD
 8 BEGIN: DOWNLOAD
 9 BEGIN: POSTDOWNLOAD
10 checking geckodriver versions:
11 BEGIN: PREDOWNLOAD
12 BEGIN: DOWNLOAD
13 BEGIN: POSTDOWNLOAD
14 checking phantomjs versions:
15 BEGIN: PREDOWNLOAD
16 BEGIN: DOWNLOAD
17 BEGIN: POSTDOWNLOAD
18 [1] "Connecting to remote server"
19 $acceptInsecureCerts
20 [1] FALSE
21
22 $browserName
23 [1] "chrome"
24
25 $browserVersion
26 [1] "79.0.3945.117"
27
28 $chrome
29 $chrome$chromedriverVersion
30 [1] "78.0.3904.11 (eaaae9de6b8999773fa33f92ce1e1bbe294437cf-refs/branch-heads/3904
      @{#86})"
31
32 $chrome$userDataDir
33 [1] "C:\\Users\\myles\\AppData\\Local\\Temp\\scoped_dir1044_653722308"
34
35
36 $`goog:chromeOptions`
37 $`goog:chromeOptions`$debuggerAddress
38 [1] "localhost:64162"
39
40
41 $networkConnectionEnabled
42 [1] FALSE
43
44 $pageLoadStrategy
45 [1] "normal"
46
47 $platformName
48 [1] "windows nt"
49
50 $proxy
51 named list()
52
53 $setWindowRect
54 [1] TRUE
55
56 $strictFileInteractability
57 [1] FALSE
58
59 $timeouts
60 $timeouts$implicit
```

```
61  [1] 0
62
63  $timeouts$pageLoad
64  [1] 300000
65
66  $timeouts$script
67  [1] 30000
68
69
70  $unhandledPromptBehavior
71  [1] "dismiss and notify"
72
73  $webdriver.remote.sessionid
74  [1] "2d4acfa5e63494664d631810c4364342"
75
76  $id
77  [1] "2d4acfa5e63494664d631810c4364342"
```

Listing 19: Fixing the error message by using one of the suggested version numbers.

You should get the output above, and a new Google Chrome window should open. It should be blank, and it should look like this:



If you get a new window that says "Chrome is being controlled by automated test software", you did this correct, and you are ready to go. If you still get error messages, it means that you either again put in the wrong chrome version (so try a different one on the suggested list!), or you have orphan R sessions. In the latter case, try closing R Studio completely and starting all over again (with the correct version number). If it still does not work, try to restart the entire compute. If it still does not work.....well, you're going to need to figure that one out for yourself!

Now that we have a new chrome browser open, we want to navigate to a website. What we will do is scrape all of the news article titles and links from a given day in the year from the Wall Street Journal. Interesting enough, WSJ has a neat archive portion of their website that organizes this info very neatly. Granted, we can take a static approach for this as well, but we will illustrate the scrapping here with the dynamic approach. We will grab the news archive from March 6th, 2011 (note how we do this in the url):

```
1 > remDr<-rd[["client"]]
2 > remDr$navigate("http://www.wsj.com/news/archive/20110306")
```

Listing 20: .

First, the first line is grabbing what is called a "handler" to the browser. This basically is like grabbing a line of communication to the web browser. We store the handler in a variable called `remDr`. Next, we use the function in the handler to navigate to the WSJ archive for March 6th, 2011. Doing so, we also notice that our web browser went to that page:



Put differently, we just controlled our web browser from R, no differently than if we were to put that url in the address bar itself and press Enter. The same action was taken from our code in R. There are other things we can do as well. In order to know where to pull the information we need, we need to now manually inspect this page so we can gain an idea of the HTML structure. If you scroll through the page in the browser, you will see that all article published on that day appear in one page (very nice indeed!). We want to pull a list of the article category, title, and link, and store this info in a file. To do so, right click on the first article and hit "inspect". You will see the following:

Notice in the HTML code that everything is well organized. Each article's information is entirely under an "article" tag. This makes life easy. All we need to do is get all article tags, iterate through each one, and extract the relevant information. Let us further explore inside a single article tag on the HTML tree to determine where the category, title, and link lie in the code. First, the category. After right clicking on the category and hitting "inspect", and some moving of the mouse back and forth until we get to the right area in the code, we can see where the category text lies in a `div` element tag as text, where the class name is "div.WSJTheme–flashline–wMh3gaMx.". While this class name is not unique in the entire HTML, it is unique in the children elements of the article tag (you can check this in a quick text search):



Next, we want to see where the title text is. We notice that upon some additional hovering in

the inspection tool in Chrome, that the title is element text in an "a" tag. This can be uniquely identified by first finding the h3 tag (since it appears there are no other h3 tags as a descendent of the article element) and then finding the only child of h3, which is the "a" tag we would need. Interestingly enough,the hyperlink to the article is also on this tag, but as a property of the "href" property:



To test our theory here, let us grab the page source, and navigate the HTML in the browser. We will not use rvest for this, but rather RSelenium's functionality (which is slightly different, and closer to that of Javascripts):

```
1  > #Get All Elements in Tree that is an Article Tag:
2  > articles<-remDr$findElements(using="tag","article")
3  >
4  > #Test the first article only:
5  > article<-articles[[1]]
6  >
7  > #See if we can find the div tag with the category text:
8  > category<-article$findChildElement(using="class","WSJTheme--flashline--wMh3gaMx"
     )
9  >
10 > #Get The tag name
11 > category$getElementTagName()
12 [[1]]
13 [1] "div"
14
15 >
16 > #We see it's a div!  Lets get the inner text:
17 > cat_text<-category$getElementText()[[1]]
18 >
19 > #Now, let us search for ONLY the h3 tag only under this article tag
20 > title<-article$findChildElement(using="tag","h3")
21 >
22 > #If we did find the h3 tag, then the a tag should be directly under it:
```

```
23 > title_link<-title$findChildElement(using="tag","a")
24 >
25 > #Grab the title of the article:
26 > title_text<-title_link$getElementText()[[1]]
27 >
28 > #Last, grab the hyperlink to the article:
29 > article_link<-title_link$getElementAttribute("href")[[1]]
30 >
31 > #Check and See!
32 > cat_text
33 [1] "WEALTH MANAGER Q&A"
34 > title_text
35 [1] "Stonehage's Armist Sees Green in Timber"
36 > article_link
37 [1] "https://www.wsj.com/articles/SB10001424052748703580004576180091880416466"
```

Listing 21: Testing an extraction method of the category and title and article link information.

Now we want to do exact this info for all the articles on the page. Let's do that!

```
1 > articles<-remDr$findElements(using="tag","article")
2 > the_data<-c()
3 >
4 > for(article in articles){
5 +     category<-article$findChildElement(using="class","WSJTheme--flashline--
        wMh3gaMx")
6 +     cat_text<-category$getElementText()[[1]]
7 +     title<-article$findChildElement(using="tag","h3")
8 +     title_link<-title$findChildElement(using="tag","a")
9 +     title_text<-title_link$getElementText()[[1]]
10 +     article_link<-title_link$getElementAttribute("href")[[1]]
11 +     the_data<-c(the_data,cat_text,title_text,article_link)
12 + }
13 >
14 > the_data<-matrix(the_data,byrow=TRUE,ncol=3)
15 > the_data<-data.frame(the_data)
16 > names(the_data)<-c("category","title","link")
17 > head(the_data)
18                   category                                            title
19 1       WEALTH MANAGER Q&A          Stonehage's Armist Sees Green in Timber
20 2                  FASHION  Akris Captures the Season With Wedding Tower Views
21 3 N.Y. HOUSE OF THE DAY                            Turn-of-the-Century Estate
22 4                 BUSINESS                       HSBC Says It Prefers London
23 5                  OPINION                        Obama's Libyan Abdication
24 6        MANAGING IN ASIA            Robson Brings Game to Thai Soccer
25                                                                          link
26 1 https://www.wsj.com/articles/SB10001424052748703580004576180091880416466
27 2 https://www.wsj.com/articles/SB10001424052748704504404576184701863626890
28 3 https://www.wsj.com/articles/SB10001424052748704570904576180961976588844
29 4 https://www.wsj.com/articles/SB10001424052748704504404576184541442696356
30 5 https://www.wsj.com/articles/SB10001424052748704005404576176861610325024
31 6 https://www.wsj.com/articles/SB10001424052748704288304576171371317870688
```

Listing 22: Doing the extraction for all articles and organizing the results in a data frame.

Last, in RSelenium, as we have mentioned, you can interact with the interface of the web browser directly. For example, we may want to navigate to one of the articles. We can certainly

do this by just copy and pasting the hyperlink and using the "navigate" function. The other way is to tell the browser to "click" on a link, so long as we "move" to the element in the HTML tree via R:

```
1  #Move to the element link on the second article, then click it:
2  articles<-remDr$findElements(using="tag","article")
3  article<-articles[[2]]
4  title<-article$findChildElement(using="tag","h3")
5  title_link<-title$findChildElement(using="tag","a")
6  title_link$clickElement()
```

Listing 23: Showing off the interactive ability of RSelenium by having the browser click on an article link.

If you go to the browser, you will now see the page for the second link up. To test other ways to interact with the browser, let us try to write code so that the browser will navigate to the main WSJ site, click on the "Search" button", open up a search box, type in "Apple and Microsoft", and then click "Search". The class names come from a series of "inspects" on each element we want to take an action on:

```
1  > remDr$navigate("http://wsj.com")
2  > search_button<-remDr$findElement(using="class","style--search-button--1U43LWwP")
3  > search_button$clickElement()
4  > text_box<-remDr$findElement(using="class","style--wsj-search-input--R7oQjW8f")
5  > text_box$clickElement()
6  > text_box$sendKeysToElement(list("Apple and Microsoft"))
7  > send_button<-remDr$findElement("class","style--search-submit--36AD7RDg")
8  > send_button$clickElement()
9  > article$clickElement()
```

Listing 24: Using RSelenium to navigate to WSJ and subsequently search for terms automatically.

While this method of web scrapping may be slower, it is more reliable. The reason is simple: we are able to "fool" the ISP and the Website into thinking it is an actual person browsing the website. This allows us to essentially see, code wise, any data that we can see directly in our browser. Dynamic web browsing is an art that takes a lot of practice. My recommendation would be to continue trying exercises found on the web for "RSelenium" examples.

**Exercises**

10. Go to http://weather.com. Using RSelenium, write code that will take the following sequence:

    - Navigate to the main website
    - Load "Wayne, NJ" weather.
    - Download the table for the 5-Day Forecast

    ONLY use RSelenium through a series of finds and clicks. You will obviously need to also use the Google Chrome's "inspect" tool.