

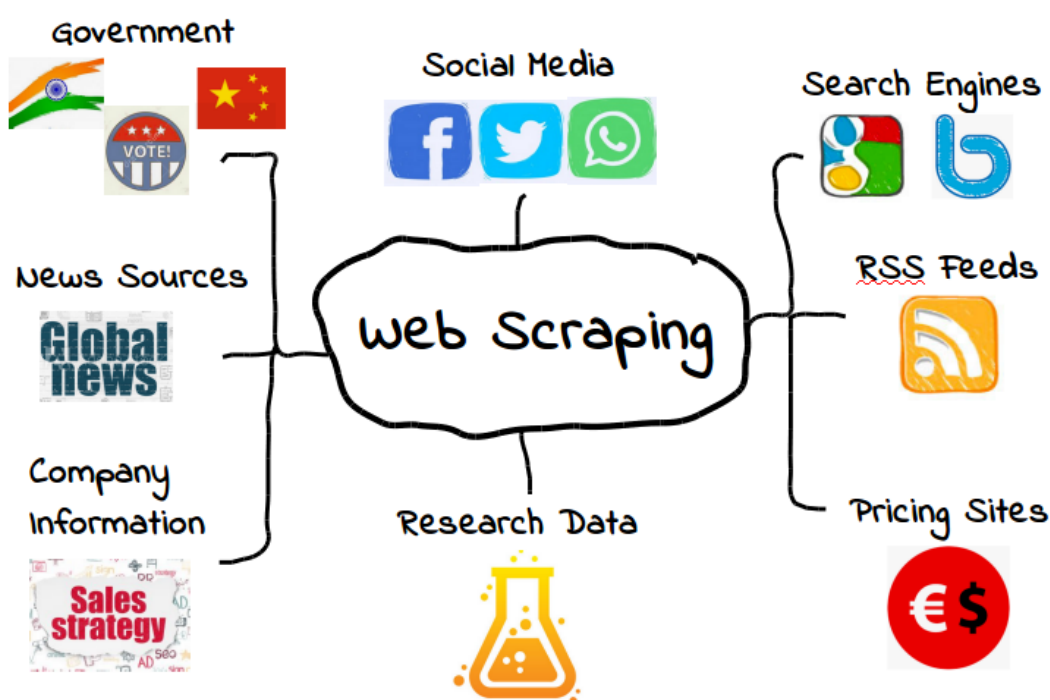
## Applied R Code

Data Science for Immediate Application

## Web Scraping in R

### Content

- [What is Web Scraping?](#)
- [Web Page Content – The basics](#)
- [Web Page Inspection](#)
- [Available R Packages](#)
- [Web Scraping Examples](#)

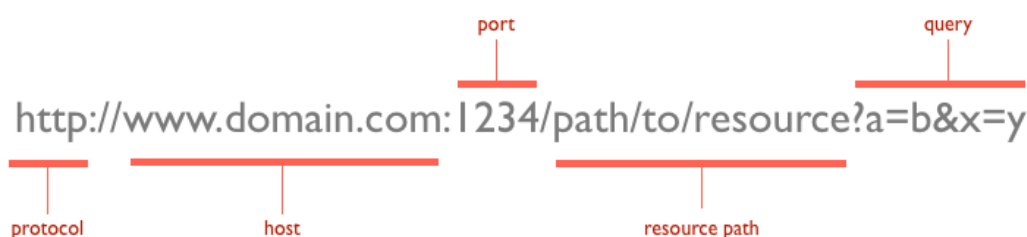


The world-wide web presents enormous amounts of data. Unfortunately, the majority of the data is not directly available for download. In response, web scraping exploits indirect means to harvest data from websites. In practice, web scraping is not unique and is totally legal. For example, web browsers rely on the Hypertext Transfer Protocol (HTTP) to fetch data and so does

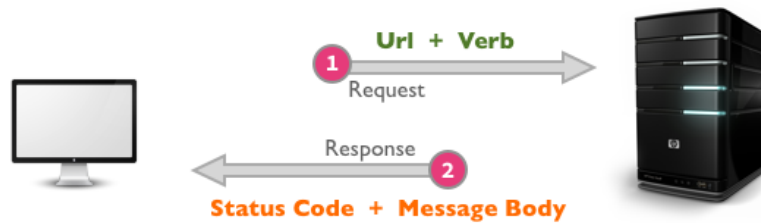
web scraping. The difference with web scraping is that the user retrieves, selects and extracts website content and data intended for browser display. This article shows how web scraping works and presents tools available in the R programming language for both manual and automated web-scraping.

### What is Web Scraping?

Web scraping involves getting a web page and extracting data from it. Specifically, a message is sent to fetch a web page using Uniform Resource Locators (URLs). A basic URL structure appears below:<sup>1</sup>



The request components include the protocol, which is typically http or https for secure communications. The host is a target website. The default port is 80, but one can be set explicitly, as shown here. The resource path specifies the server path to the data and the query is the data request action or verb.



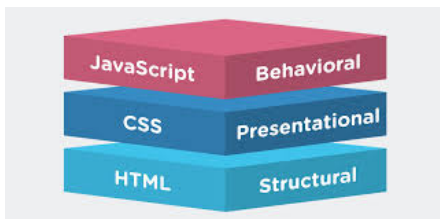
A web page request is a simple component of web scrapping. Next, a server response will deliver a status message with web page content. The web page content must be searched, be it manually or automatically. Hence, web page crawling is a key feature of web scrapping. Next, the web page content will be parsed, extracted and reformatted.

Web scrapping, for example, might first retrieve a web page and then extract contact names and phone numbers. In another instance, the focus might be to grab a research data table or a collection of tables. Finally, another effort might seek to extract intel using text strings found in multiple news or social media reports.

### ***Web Page Content – The basics***

The response to a URL request is the web page delivered in a

HTML document message. HTML stands for HyperText Markup Language and defines the content and structure of a webpage.



“Hyper Text” in HTML refers to hyperlinks that connect webpages to one another, either within a single website or between websites, to populate page content. Meanwhile, “markup” text includes special page elements or “CSS selectors” such as <head>, <title>, <section>, <body>, <div>, <table>, <p>, <img> and may others. Web crawling and scrapping involves finding and extracting from these elements as needed.

Other language technologies also populate the web page besides HTML. The technologies include CSS, which describe a webpage’s presentation appearance and JavaScript for web page functionality. An analyst inspects all web page content using a web browser tool to find what parts will be scraped.

### ***Web Page Inspection***

Inspect web page content using web browser tools to find what parts will be scraped. For example, the following web site image (<https://oilprice.com/oil-price-charts>) shows petroleum product prices from around the world. The lower half of the image shows the web site HTML content in your web browser inspector (click to enlarge):

**Oil Price Charts**

**Futures & Indexes**

Opec Members (Daily Pricing)

International Prices

Canadian Blends

United States Blends

Opec Members

| Futures & Indexes       | Last  | Change | % Change | Last Updated     |
|-------------------------|-------|--------|----------|------------------|
| WTI Crude               | 55.98 | +1.19  | +2.17%   | (16 Hours Delay) |
| Brent Crude             | 66.25 | +1.68  | +2.6%    | (16 Hours Delay) |
| Mars US                 | 63.19 | +1.63  | +2.65%   | (16 Hours Delay) |
| Opec Basket             | 64.28 | +1.34  | +2.13%   | (2 Days Delay)   |
| Canadian Crude Index    | 42.29 | +0.00  | +0.00%   | (15 Hours Delay) |
| DME Oman                | 66.73 | +1.16  | +1.77%   | (16 Hours Delay) |
| Urals                   | 63.07 | +0.82  | +1.32%   | (1 Day Delay)    |
| Mexican Basket          | 57.39 | +0.81  | +1.43%   | (2 Days Delay)   |
| Indian Basket           | 65.04 | +1.62  | +2.55%   | (2 Days Delay)   |
| Western Canadian Select | 43.76 | +0.11  | +0.25%   | (17 Hours Delay) |
| Dubai                   | 64.05 | +0.76  | +1.20%   | (2 Days Delay)   |
| Brent Weighted Average  | 64.19 | +0.76  | +1.20%   | (2 Days Delay)   |
| Louisiana Light         | 63.60 | +1.41  | +2.27%   | (2 Days Delay)   |
| Coastal Grade A         | 45.00 | +1.00  | +2.27%   | (1 Day Delay)    |
| Domestic Swt. @ Cushing | 52.00 | +1.00  | +1.96%   | (1 Day Delay)    |
| Giddings                | 45.75 | +1.00  | +2.23%   | (1 Day Delay)    |
| ANS West Coast          | 64.27 | +1.20  | +1.90%   | (2 day delay)    |
| Gulf Coast HSFO         | 62.77 | +0.53  | +0.85%   | (2 Days Delay)   |
| Natural Gas             | 2.656 | +0.046 | +1.76%   | (16 Hours Delay) |
| Heating Oil             | 2.017 | +0.049 | +2.51%   | (16 Hours Delay) |
| Gasoline                | 1.738 | +0.054 | +3.23%   | (16 Hours Delay) |
| Ethanol                 | 1.329 | +0.003 | +0.23%   | (17 Hours Delay) |

09:32 am CST 16/02/2019

**Heating Oil**

**2.017 +2.51%**

1D | 1WK | 1M | 1YR | Max

View: Brent Crude Reset chart

Add a compar

**Inspector** Console Debugger {} Style Editor Performance Memory Network Storage Accessibility

Search HTML

Filter Styles

```

<div class="oilprices_centercolumn">
  <table class="oilprices_table" data-id="1">
    <tbody>
      <tr class="row_holder">
        <tr class="stripe show_graph update_on_load" data-id="45" data-spreadsheet="Crude Oil WTI">
          <td class="flag_holder">
            <td class="last_price" data-price="55.98">55.98</td>
            <td class="change_up flat_change_cell" dir="ltr">+1.19</td>
            <td class="change_up percent percent_change_cell" dir="ltr">+2.17</td>
            <td class="last_updated" data-stamp="1550271780">(16 Hours Delay)</td>
          <tr class="stripe show_graph update_on_load" data-id="46" data-spreadsheet="Crude Oil Brent">
            <tr class="stripe show_graph update_on_load" data-id="50" data-spreadsheet="Gulf Sour Crude Oil">
            <tr class="stripe show_graph update_on_load" data-id="29" data-spreadsheet="Opec Basket Price">
            <tr class="stripe show_graph update_on_load" data-id="68" data-spreadsheet="Canadian Crude Index">

```

Failure with checking formfill: TypeError: b is null ln: 71

The inspector is launched in FireFox by right-clicking the data table and choosing “Inspect Element.” Other web browsers have selections to “View Page Source.” The next image shows the page source for the first table (click to enlarge):

If you don’t know HTML, this may look daunting! But look close. All the table text and data is there to be grabbed. Meanwhile, another essential tool for web page inspection is [SelectorGadget](#). SelectorGadget is an open source tool that makes it easy to identify and define CSS content. Just install the browser extension as instructed. The tool simplifies inspection of complicated web sites to just a few point-and-click actions and it makes it a breeze to locate content to scrape.

## Available R Packages

Several R packages are available for webscraping:

| Package Name | Retrieve? | Parse? | Must Know? |
|--------------|-----------|--------|------------|
| rvest        | YES       | YES    | YES        |
| RCurl        | YES       | NO     | YES        |
| XML          | Limited   | YES    | YES        |
| rjson        | NO        | YES    | YES        |
| RJSONIO      | NO        | YES    | Optional   |
| httr         | YES       | YES    | Optional   |
| xml2         | NO        | YES    | Optional   |
| selectr      | NO        | YES    | Optional   |

`rvest` is a set of wrapper functions around the `xml2` and `httr` packages. The rest of the tutorial focuses on `rvest` exclusively. The main functions are:

- `read_html()`: read a webpage into R as XML (document and nodes)
- `html_nodes()`: extract pieces out of HTML documents using XPath and/or CSS selectors
- `html_attr()`: extract attributes from HTML, such as `href`
- `html_text()`: extract text content

## Web Scraping Examples

The first step is to load the `rvest` package and to read the target web page:

```

1 library(rvest)
2
3 url <- "https://oilprice.com/oil-price-charts"
4 webpage <- read_html(url)
5
6 webpage
7 {xml_document}
8 <html lang="en">
9 [1] <head>\n<script>\n  (function(i,s,o,g,r,a,m){i['GoogleAnalyticsObject']=r;i[r]=i[r]||fun ...
10 [2] <body id="pagetop" class="oilprices loggedout">\n<!-- Google Tag Manager (noscript) -->\n ...

```

The first two lines above are the first HTML lines of the target web page. `rvest` has read or captured the entire HTML content of the target webpage. Next, the webpage nodes are read to extract to find data table content...the product names and prices. The `rvest` command `html_nodes()` is told to seek specific nodes as defined by CSS selectors used in the function calls. The selector string used was defined by SelectorGadget after clicking on a price or name element on the web page. The remaining command `html_text()` converts node content to text and other commands use basic R syntax to reshape the data for use in R:

```

1 value <- webpage %>%
2   html_nodes(css = ".last_price") %>%
3   html_text()
4
5 name <- webpage %>%
6   html_nodes(css = "td:nth-child(2)") %>%
7   html_text() %>%
8   .[c(41:72, 74, 76:77, 79, 81:82, 84:85, 87:89, 91, 94:96, 98:99, 101:103, 105:107, 109:112, 114:116, 118:120, 122
9     126:128, 130:131, 133, 135:142, 145:160, 162:164, 166:167, 169:170, 172, 174,
10    176:178, 180:182, 184:186, 188:189, 191, 193:196, 198:199, 201:213)]
11
12 prices <- tibble(name = name,
13                 last_price = value)
14 prices
15 # A tibble: 138 x 2
16   name                last_price
17   <chr>              <chr>
18 1 WTI Crude          55.98

```

|    |                            |       |
|----|----------------------------|-------|
| 19 | 2 Brent Crude              | 66.25 |
| 20 | 3 Mars US                  | 63.19 |
| 21 | 4 Opec Basket              | 64.28 |
| 22 | 5 Canadian Crude Index     | 42.29 |
| 23 | 6 DME Oman                 | 66.73 |
| 24 | 7 Urals                    | 63.07 |
| 25 | 8 Mexican Basket           | 57.39 |
| 26 | 9 Indian Basket            | 65.04 |
| 27 | 10 Western Canadian Select | 43.76 |
| 28 | 11 Dubai                   | 64.05 |
| 29 | 12 Brent Weighted Average  | 64.19 |
| 30 | 13 Louisiana Light         | 63.60 |
| 31 | 14 Coastal Grade A         | 45.00 |
| 32 | 15 Domestic Swt. @ Cushing | 52.00 |

The next example uses the pipe operator to link node parsing, data extracting and formatting. In this case, the extraction is on the 13th data table returned by the `html_table()` function. The table scraping extracts the spot price and price change data for US crude oil blends:

```

1 tables13 <- webpage %>%
2   html_nodes("table") %>%
3   html_table(fill = TRUE) %>%
4   .[[13]] %>%
5   .[c(3:18, 21:23, 26:27, 30:31, 34, 37, 40:42, 45:47, 50:52, 55:56, 59, 62:65, 68:69, 72),] %>%
6   `colnames<-`(c("index", "name", "last", "change", "percent", "updated")) %>%
7   as.tibble() %>%
8   separate(col = "percent", into = c("percent", "update"), sep = "\\(") %>%
9   select(name, last, change, percent) %>%
10  mutate(last = as.numeric(last)) %>%
11  mutate(change = as.numeric(change)) %>%
12  mutate(percent = as.numeric(substr(percent,1,5)))
13
14 table13
15 A tibble: 44 x 4
16   name                last change percent
17   <chr>                <dbl> <dbl> <dbl>
18 1 West Texas Sour      49.54  1.18  2.44
19 2 West Texas Intermediate 52.04  1.18  2.320
20 3 Upper Texas Gulf Coast 39.84  1.18  3.05
21 4 Texas Gulf Coast Light 50.54  1.18  2.39
22 5 South Texas Sour      45.93  1.18  2.64
23 6 North Texas Sweet     49.5   0.25  0.51
24 7 North Texas Sour      39.91  0.87  2.23
25 8 Eagle Ford Pipeline    52.04  1.18  2.320
26 9 Eagle Ford Condensate  51.04  1.18  2.37
27 10 Eagle Ford           53.49  1.18  2.260
28 11 Tx. Upper Gulf Coast  45.75  1     2.23
29 12 South Texas Light     45.75  1     2.23
30 13 W. Tx./N. Mex. Inter.  52     1     1.96
31 14 South Texas Heavy     45.5   1     2.25
32 15 W. Cen. Tx. Inter.    52     1     1.96
33 16 East Texas Sweet      49.25  1     2.070
34 17 Arkansas Sweet        47.75  0.25  0.53
35 18 Arkansas Sour         46.75  0.25  0.54
36 19 Arkansas Ex Heavy     42.75  0.25  0.59
37 20 Buena Vista           65.21  1.11  1.73
38 21 Midway-Sunset         60.53  1.11  1.87
39 22 Williston Sweet       45.75  0.75  1.67
40 23 Williston Sour        41.86  0.75  1.82
41 24 Utah Black Wax        39.11  0.51  1.32
42 25 Four Corners          49.68  0.5   1.02
43 26 Colorado D-J Basin    50     0.75  1.52
44 27 Colorado South East   41.25  0.25  0.61
45 28 Colorado West         47.41  0.51  1.09
46 29 NW Kansas Sweet       42.25  0.25  0.6
47 30 SW Kansas Sweet       42.75  0.25  0.59
48 31 South Central Kansas  46.5   2     4.49
49 32 Delhi/N. Louisiana    49     1     2.08
50 33 South Louisiana       50.5   1     2.02
51 34 North Louisiana Sweet  47.75  0.25  0.53
52 35 Michigan Sour         44     1     2.33

```

|    |    |                       |       |      |       |
|----|----|-----------------------|-------|------|-------|
| 53 | 36 | Michigan Sweet        | 48.75 | 1    | 2.09  |
| 54 | 37 | Nebraska Sweet        | 46.51 | 0.51 | 1.11  |
| 55 | 38 | Oklahoma Sweet        | 52.04 | 1.18 | 2.320 |
| 56 | 39 | Oklahoma Sour         | 37.5  | 0.25 | 0.67  |
| 57 | 40 | Western Oklahoma Swt. | 51.25 | 1    | 1.99  |
| 58 | 41 | Oklahoma Intermediate | 51    | 1    | 2     |
| 59 | 42 | Wyoming General Sour  | 42.1  | 0.75 | 1.81  |
| 60 | 43 | Wyoming General Sweet | 48.75 | 0.75 | 1.56  |
| 61 | 44 | Central Montana       | 46.98 | 0.51 | 1.1   |

1. Image from <https://code.tutsplus.com/tutorials/http-the-protocol-every-web-developer-must-know-part-1-net-31177> 

---

## Applied R Code

Proudly powered by WordPress.