

Five easy steps for *scraping data* from web pages.

As published in Benchmarks RSS Matters, November 2013

<http://web3.unt.edu/benchmarks/issues/2013/11/rss-matters>

Jon Starkweather, PhD

Jon Starkweather, PhD
jonathan.starkweather@unt.edu
Consultant
Research and Statistical Support



<http://www.unt.edu>

RSS
Research and Statistical Support

<http://www.unt.edu/rss>

RSS hosts a number of “Short Courses”.
A list of them is available at:
<http://www.unt.edu/rss/Instructional.htm>

Those interested in learning more about R, or how to use it, can find information here:
http://www.unt.edu/rss/class/Jon/R_SC

Five easy steps for *scraping* data from web pages.

In a perfect world, all data would be easily available to everyone as comma separated values (CSV) files. Unfortunately, Earth is not yet a perfect world. Occasionally, some interesting data is unavailable as a CSV download, but is available and / or displayed on a web page. The term scraping data refers to the process of parsing the HTML source code of an available web page in order to extract, retrieve, or scrape some specific data from the web page. In truth, the title of this article is a bit misleading. The functions used to scrape data are fairly straightforward and easy to use; but, there is a significant assumption when using them. The time consuming assumption is you need to know where the desired data is located among the lines of HTML code of the web page you are scraping. However, once the line, or lines, have been identified; the R script used to scrape the data and put it into a manageable R format is very easy to use and can be reused over time. The example below illustrates the process of scraping the DOW Jones Industrial average from the Reuters Commodities (2013) website; specifically, the following web page: <http://www.reuters.com/finance/commodities/energy#oil> (displayed below). Keep in mind; scraping data does not require the target web page to be open in a browser. In fact, the script in this article, and other scripts like it, will work without ever opening a browser - only an internet connection is required. Of course, the target web page needs to be public (i.e. not behind a log in or encrypted).

The screenshot shows the Reuters Commodities website. It features several data tables and a market indices section. A green box highlights the DOW index value of 14,805.05 in the Market Indices table.

| Commodity | Exch | Currency | Expire | Last | Trade Date/Time | Net Chg | Open | High | Low |
|---------------------------|------|----------|--------|--------|-----------------|---------|--------|--------|--------|
| LIGHT CRUDE COM1 Sep13 | NYM | USD | 09/20 | 108.73 | 08/27 14:20 | +2.81 | 106.14 | 109.32 | 105.88 |
| NO 2 HT OIL COM1 Aug13 | NYM | USD | 08/30 | 3.15 | 08/27 14:20 | +0.07 | 3.09 | 3.16 | 3.08 |
| BRENT CRUDE JUL4 Aug13 | IEU | USD | 08/15 | 113.92 | 08/27 14:20 | +3.19 | 111.16 | 114.35 | 110.73 |
| GAS OIL APR1 Sep13 | IEU | USD | 09/12 | 960.50 | 08/27 14:20 | +20.25 | 941.50 | 963.00 | 937.75 |
| GASOLINE JAIN Aug13 | TCE | JPY | 08/23 | 77,420 | 08/28 13:28 | +670 | 76,610 | 77,420 | 76,350 |
| KEROSENE JAIN Aug13 | TCE | JPY | 08/23 | 77,500 | 08/28 07:18 | +230 | 76,900 | 77,500 | 76,900 |

Data as of 2:31pm EDT (Delayed at least 20 minutes).

| Commodity | Exch | Currency | Expire | Last | Trade Date/Time | Net Chg | Open | High | Low |
|---------------------------|------|----------|--------|------|-----------------|---------|------|------|------|
| NATURAL GAS COM1 Aug13 | NYM | USD | 08/28 | 3.53 | 08/27 14:20 | +0.01 | 3.51 | 3.53 | 3.45 |

Data as of 2:31pm EDT (Delayed at least 20 minutes).

| Commodity | Exch | Currency | Expire | Last | Trade Date/Time | Net Chg | Open | High | Low |
|--------------------------------|------|----------|--------|-------|-----------------|---------|------|------|------|
| PJM ELECTRICITY NYMEX Oct12 | USYF | USD | 10/30 | 39.40 | 10/08 20:00 | +0.00 | 0.00 | 0.00 | 0.00 |

Data as of 2:31pm EDT (Delayed at least 20 minutes).

MARKETS

U.S. EUROPE ASIA SECTORS

Market Indices

| | | | |
|-------------|-----------|---------|--------|
| DOW | 14,805.05 | -141.41 | -0.95% |
| S&P 500 | 1,634.23 | -22.55 | -1.36% |
| NASDAQ | 3,585.38 | -72.19 | -1.97% |
| TR US INDEX | 149.26 | -1.93 | -1.28% |

Currencies

| | | |
|---------|--------|--------|
| EUR/USD | 1.3391 | +0.18% |
| GBP/USD | 1.5535 | -0.25% |
| USD/JPY | 97.070 | -1.45% |

Commodities

| | | | |
|------|----------|--------|--------|
| GOLD | 1,420.10 | +27.60 | +1.98% |
| OIL | 108.73 | +2.81 | +2.65% |
| CORN | 500.25 | -15.50 | -3.01% |

Market News

The specific data retrieved in this example will be the DOW, which is listed in the small Markets table


```
dow.df <- data.frame(matrix(rep(NA,3), ncol = 3))
names(dow.df) <- c("string.date","numeric.date","DOW")
dow.df
  string.date  numeric.date  DOW
1          NA             NA   NA
```

Below, we retrieve the date and time using the 'Sys.time' function and make sure to store the numeric version as well as the character string version.

```
dow.df[,1] <- Sys.time()
dow.df[,2] <- as.numeric(Sys.time())
```

Finally, we can convert our data into numeric and put it into the data frame we created.

```
dow.df[,3] <- as.numeric(new.line.3)
dow.df
  string.date  numeric.date  DOW
1 2013-08-27 14:58:10    1377633491 14805.05
```

Step 4: System time and a data frame.

We can then save or export the data by setting the working directory to the location we want to store the file and using the 'write.table' function.

```
setwd("C:/Users/jds0282/Desktop/")
write.table(dow.df, file = "dow.df.txt", sep = ",", na = "NA",
  dec = ".", row.names = TRUE, col.names = TRUE)
```

Conclusions

Keep in mind, although it may take significant effort to identify the line number of the data of interest, once the script has been written and checked, it can be used repeatedly (e.g. each day) to retrieve the data of interest (i.e. to build a time series data file). The only real problem which can occur is when the HTML source code is changed, in other words, if the web page author(s) update(s) the layout of the page. Then, of course, it would be necessary to verify the line number of the data of interest and re-check the script to make sure it returns the desired information.

As stated above, there are other ways of accomplishing what was accomplished in this article; the 'RCurl' package is apparently quite popular. All of the functions used in this article are available with a base install of R - the functions are available in the 'base' package. For more information on what R can do, please visit the Research and Statistical Support Do-It-Yourself Introduction to R² course website. An Adobe.pdf version of this article can be found here³.

²http://www.unt.edu/rss/class/Jon/R_SC/

³<http://www.unt.edu/rss/rssmattersindex.htm>

Until next time; *“information wants to be free”*...

References & Resources

Lang, D. T. (2013). Package RCurl. Package documentation available at: <http://cran.r-project.org/web/packages/RCurl/index.html> and further information available at: <http://www.omegahat.org/RCurl/>

ProgrammingR.com (2013) Webscraping using readLines and RCurl. Available at: <http://www.programmingr.com/content/webscraping-using-readlines-and-rcurl/>

Reuters. (2013). Commodities: Energy. Retrieved on August 27, 2013 from: <http://www.reuters.com/f>

This article was last updated on November 5, 2013.

This document was created using \LaTeX