

# Introduction to Webscraping with R

By Robert Vesco

>> *For Access to R Code, Please Open this Presentation  
in Dedicated PDF Application and Click on Pin <<*





# Outline

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- 1 Why Use R for Webscraping
- 2 Why XML, XPath Approach
- 3 The Basics of Webscraping
- 4 R Example
- 5 RCurl
- 6 Practical Advice
- 7 References



# Why Use R for Webscraping

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- No external languages or scripts needed
- Makes workflow more efficient
- Easier to share with othe[R] colleagues
- Can accomplish most webscraping needs quickly and efficiently



# Outline

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- 1 Why Use R for Webscraping
- 2 Why XML, XPath Approach
- 3 The Basics of Webscraping
- 4 R Example
- 5 RCurl
- 6 Practical Advice
- 7 References



# Why XML, XPATH Approach

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPATH  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- Faster than using regular expressions
- More robust
- Nearly all languages now support XPATH approach
- HTML code in the wild getting better all the time – and hence makes XPATH more reliable
- Can and should be used with regular expressions



# Outline

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- 1 Why Use R for Webscraping
- 2 Why XML, XPath Approach
- 3 The Basics of Webscraping**
- 4 R Example
- 5 RCurl
- 6 Practical Advice
- 7 References



# HTML - The Code Behind the Web

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References



**Robert V**

**Joined:** April 12, 2010

"1) I'm a love[r] not a hate[r] 2)

```
1 <li id="member_4403063">
2   <div id="image_4403063" class="profphoto">
3     <a href="http://www.meetup.com/R-users-DC/members/4403063/">
4       
7     </div>
8     <div class="memberItem">
9       <div class="memberInfo">
10        <div id="memberInfo_4403063" class="D_title">
11          <a class="memName" href="http://www.meetup.com/R-users-DC/members/4403063/↵
12            ↵>Robert V</a>
13          <span id="role_4403063" class="memRole">
14            </span>
15          <ul class="D_less memStats">
16            <li>
17              <span class="bold">Joined</span>: April 12, 2010
18            </li>
19            <p class="D_less">
20              "1) I'm a love[r] not a hate[r]
21              2)
```



# Using XPATH to Select Items

We want to select just my name

```
1 .//*[@class='memName ']
```

or use fuller path

```
1 .//*div/div/div/ etc .... /a[@class='memName ']
```

Will both pull out "Robert V" from this code:

```
1 <a class="memName" href="http://www.meetup.com/R-users ↵  
  -DC/members/4403063/">Robert V</a>
```

BUT - more importantly, the above code will also pull out every name if you're looking at all the code on the page!

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References





# XPATH Tutorials

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- <http://www.zvon.org/xxl/XPathTutorial/General/examples.html>
- <http://www.w3schools.com/xpath/>



# Important XML Functions in R

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

`htmlTreeParse()` parses file, cleans malformed HTML, and make it available for querying

```
1 web <- htmlTreeParse(chrFileUrl, error=function(...){}, useInternalNodes =  
TRUE, trim=TRUE)
```

`getNodeSet()` takes xml documents and queries it for particular nodes (html tags)

```
1 nArticles <- getNodeSet(web, "//*[div[@class='D_title']]")
```

`xmlValue()` takes nodes from `getNodeSet` and extracts text value

```
1 xmlValue(nArticles[[1]])
```

`xmlGetAttr()` get attributes from html tags such as href links or class names

```
1 xmlGetAttr(nArticles[[1]], "href")
```

Lastly, there exists convenience functions such as `xpathSApply()`



# Some Tools - Firebug + FireXPath

- Allows you to select items on a webpage and inspect their underlying tags. Also allows you to query your XPATH to see what it will select!

The screenshot shows a web browser displaying a Meetup profile for Robert V. The profile includes a name, a photo, a bio, and a list of other members. The Firebug extension is open, showing the DOM tree and the XPATH query used to select the member's bio text.

Washington, DC

105 R Users

2 comments

Meetups Founded August 10, 2009

Membership dues \$0.00 - on joining (more info)

Robert V

Joined: April 12, 2010

I'm a love[r] not a hate[r] 2) I use for R for nearly all my webscraping and data analysis needs. "

Consola HTML CSS Script DOM Net XPath

Top Window XPath: `//*[@id=memberInfo_4403063]`

```
<div class="memberiten">
  <div class="memberInfo">
    <div id="memberInfo_4403063" class="D_title">
      <a class="memname" href="http://www.meetup.com/R-users-DC/members/4403063/">Robert V</a>
      <span id="role_4403063" class="memRole"/>
    <div class="D_less memStats">
      <ul>
        <li>
          <span class="bold">Joined</span>
          : April 12, 2010
        </li>
      </ul>
      <p class="D_less"> "1) I'm a love[r] not a hate[r] 2) I use for R for nearly all my webscraping and data analysis needs. " </p>
    </div>
  </div>
</div>
</li>
<li id="member_12237852">
<li id="member_2906765">
<li id="member_10173172">
<li id="member_10271290">
<li id="member_10227539">
```



# Some Tools - Selector Gadget

Introduction to  
Webscraping with R

By Robert Vesco

Why Use R for  
Webscraping

Why XML,  
XPath Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- Allows you to select multiple elements on a screen
- Useful for very complicated layouts

The screenshot shows a Meetup page for a group in Washington, DC. The page title is "Members" with 105 R Users. A "Members" list is visible, including "Geite Kim" (joined June 3, 2010) and "Lynn Flowers" (joined February 27, 2010). A "Selector Gadget" window is overlaid on the page, displaying the CSS selector `//*[@contains(concat(" ", @class, " "), concat(" ", "D_metaLabel", " "))]/*/[*(@id = "member_440...)]`. The window also shows the URL `http://www.meetup.com/says:` and a message: "The CSS selector '.D\_metaLabel, #member\_4403063' as an XPath is shown below. Please report any bugs that you find with this converter." The gadget interface includes "OK" and "Cancel" buttons, and a search bar at the bottom containing the selector `.D_metaLabel, #member_4403063`.



# Outline

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- 1 Why Use R for Webscraping
- 2 Why XML, XPath Approach
- 3 The Basics of Webscraping
- 4 R Example**
- 5 RCurl
- 6 Practical Advice
- 7 References



# General Steps

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- Select and download the pages you want
- Query the document
- Select which items you want and save to dataframe
- Repeat



# Find Which Pages Have the Info You Need + Download

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

```
1  #Seq Variables
2  cStartSeq ← 0
3  cEndSeq ← 100
4  cStepSeq ← 20
5  #link variables
6  chrURLPrefix ← "http://www.meetup.com/R-users-DC/members/?offset="
7  chrURLSuffix ← "&desc=1&sort=chapter_member.atime"
8  #Files will be read and save to this path. Make sure it exists on your computer
9  #or change path to wherever you save this script to!!
10 chrSetDir = "~/R/R_Meetup/"
11 #setwd(chrSetDir) #commented out for sweave
12 chrDir ← paste(chrSetDir,"RawData/",sep="")
13 #Check to see if folder exists, else create it.
14 if("RawData" %in% dir(chrSetDir)==FALSE){dir.create(chrDir)}
15 for(w in seq(cStartSeq,cEndSeq,cStepSeq))
16 {
17   #Create URL that will download page
18   url ← paste(chrURLPrefix,w,chrURLSuffix,sep="")
19   #Create name for URL. Important because sometimes URL names have illegal ←
     characters
20   #or lengths for files systems.
21   urlName ← paste(chrDir,w,".html",sep="")
22   #Without error catching script will crash. Websites frequently time out!
23   err ← try(download.file(url,destfile = urlName,quiet = TRUE), silent = TRUE)
24   if(class(err)=="try-error")
25   {
26     #you may be hitting the server too hard, so backoff and try again later.
27     Sys.sleep(5) #in seconds, adjust as necessary
28     try(download.file(url,destfile = urlName,quiet = TRUE), silent = TRUE)
29   }
30 }
```



# Process Pages and Extract Data

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

```
1 require(XML)
2 require(xtable)
3 vFiles ← list.files(chrDir)# put files in rawdata folder into vector and get ←
   length
4 iLenFilesList ← length(vFiles) #create list to store dataframes
5 ls ← list()
6 for(i in 1:iLenFilesList)
7 {
8   #each i will pull a different URL
9   url ← vFiles[i]
10  chrFileUrl ← paste(chrDir, url, sep="")
11  # this function works on dirty html, adding closing tags and such.
12  web ← htmlTreeParse(chrFileUrl, error=function(...){}, useInternalNodes = TRUE, ←
   encoding = "UTF-8", trim=TRUE)
13  #Use vectorized function to get names
14  vNames ← xpathSApply(web, '//*[@class="memName"]', xmlValue)
15  #Same as above, but use regex to clean up a bit
16  vDates ← gsub("Joined: |\\r\\n", "", xpathSApply(web, '//*[@class="D_less ←
   memStats"]/li', xmlValue))
17  #Since not every person has a quote, we break the problem into part getting ←
   chunks of code
18  vQuote2 ← getNodeSet(web, "//*[@div[@class='D_title ']")
19  #now we look for quotes - notice ".//*" this means subquery -- IMPORTANT!
20  vQuote3 ← sapply(vQuote2, function(x) xpathSApply(x, ".//p[@class='D_less ']", ←
   xmlValue))
21  # we get list() for node with no quotes. Replace list() with NULL
22  vQuote4 ← gsub('|\\r|\\n|["']', "", sapply(vQuote3, function(x) ifelse(is.list(x), NA, ←
   x)))
23  #add df to list. This is ok for small scrapes, but for larges ones, you need to ←
   write to file or db
24  ls[[i]] ← data.frame(Name=vNames, Date=vDates, Quote=vQuote4, stringsAsFactors=←
   FALSE)
25 }
```





# Combine Data from Each Loop

```
1 #combine df
2 df ← do.call(rbind,ls)
3 #sample output for latex
4 library(Hmisc)
5 df2 ← df[1:3,]
6 latex(df2, file='', col.just=c("l","l","p{2in}"))
```

df2	Name	Date	Quote
1	JOEL ROBERTS	June 9, 2010	My name is Joel Roberts. I am a friend of Bryan Stroube who told me about the DC useR Group. I work with several systems that use XML for information exchange between dissimilar computer / software systems.
2	Arun	May 1, 2010	NA
3	Travis M	April 16, 2010	NA



# Outline

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- 1 Why Use R for Webscraping
- 2 Why XML, XPath Approach
- 3 The Basics of Webscraping
- 4 R Example
- 5 RCurl**
- 6 Practical Advice
- 7 References



# RCurl Package

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- More flexible, allows one to modify headers, referers, etc...
- Beyond http: https, ftp, sftp, scp, ldap, etc.....
- Come from c library libcurl so fast, extensive, and actively developed
- Can use persistent connections, cookies, and process requests as they come in rather than sequentially



# Some RCurl Examples

`getURL()` allows the use of https whereas built-in R functions do not (recent R versions may be different, internet2 method must have valid cert)

```
1 txt = getURL("https://www.twitter.com", ssl.verifyhost=0, ssl.verifypeer←  
=0)
```

`getCurlHandle()` handles allow persistent connections and settings to be used across repeated call to same server which is similar to pass list of arguments, but potentially with better network connectivity.

```
1 curl = getCurlHandle(cookie=cookie, useragent= "Mozilla/5.0 (Windows; U; ←  
Windows NT 5.1; en-US; rv:1.8.1.6) Gecko/20070725 Firefox/2.0.0.6")  
2 txt = getURL("http://www.meetup.com/R-users-DC/members/", curl=curl, .opts ←  
= list(verbose = TRUE))  
3 txt2 = getURL("http://www.meetup.com/R-users-DC/members/", curl=curl, .opts←  
= list(verbose = TRUE))
```

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References



# Outline

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- 1 Why Use R for Webscraping
- 2 Why XML, XPath Approach
- 3 The Basics of Webscraping
- 4 R Example
- 5 RCurl
- 6 Practical Advice**
- 7 References



# Warning and Practical Advice

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- Be Nice! Go easy on the servers – especially if you're using rCURL
- Check to make sure website does not prohibit scrapping. See if they already have an API. Else, send an email to web owners.
- If login required, then need to use rCurl package.
- Always use a proxy - if for nothing else you don't want your home address to accidently get blocked (ie google automated query blocking).



# Warning and Practical Advice

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- Consider setting up your old computer or laptop to run scraping. Frees your main computer up. Set it to email you if it crashes.....
- Get familiar with error catching and debugging
- Even if you try to make your code robust, large jobs will require several retweakings of code, because ex ante, you don't know all the possible permutations. Your code will become more generalized with more exception handling.
- XML packages have slightly different interpretations. Hence going from one language to the next may require 'slightly' different queries



# Outline

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach

The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- 1 Why Use R for Webscraping
- 2 Why XML, XPath Approach
- 3 The Basics of Webscraping
- 4 R Example
- 5 RCurl
- 6 Practical Advice
- 7** References





# References

Introduction  
to  
Webscraping  
with R

By Robert  
Vesco

Why Use R for  
Webscraping

Why XML,  
XPath  
Approach


The Basics of  
Webscraping

R Example

RCurl

Practical  
Advice

References

- R code used in this presentation, click on pin 
- **Tutorials,**  
<http://www.zvon.org/xxl/XPathTutorial/General/examples.html>  
<http://www.w3schools.com/xpath/>
- **Firebug**  
<http://getfirebug.com/>  
<https://addons.mozilla.org/en-US/firefox/addon/11900/>
- **SelectorGadget**  
<http://www.selectorgadget.com/>
- **XML package docs**  
<http://www.omegahat.org/RXML/>
- **rCurl Package**  
<http://www.omegahat.org/RCurl/>