

# MPP-E1180 Lecture 7: Web Scraping + Transforms

Christopher Gandrud

16 October 2015

# Objectives for the week

- ▶ Assignments
- ▶ Review
- ▶ Intro to web scraping
- ▶ Processing strings, including an intro to regular expressions
- ▶ Data and data set transformations with dplyr

## Assignment 2

**Proposal** for your Collaborative Research Project.

**Deadline:** 23 October

**Submit:** A (max) 2,000 word proposal created with **R Markdown**.  
The proposal will:

- ▶ Be written in R Markdown.
- ▶ State your research question. And justify why it is interesting.
- ▶ Provide a basic literature review (properly cited with BibTeX).
- ▶ Identify data sources and appropriate research methodologies for answering your question.

As always, submit the entire GitHub repo.

# Assignment 3

**Purpose:** Gather, clean, and analyse data

**Deadline:** 13 November 2015

You will submit a GitHub repo that:

- ▶ Gathers web-based data from at least **two sources**. Cleans and merges the data so that it is ready for statistical analyses.
- ▶ Conducts basic descriptive **and** inferential statistics with the data to address a relevant research question.
- ▶ Briefly describes the results including with **dynamically** generated tables and figures.
- ▶ Has a write up of **1,500 words maximum** that describes the data gathering and analysis, It also will use literate programming.

# Assignment 3

This is ideally a **good first run** at the data gathering and analysis parts of your final project.

# Review

What is open public data?

- ▶ Name one challenge and one opportunity presented by open public data.

What is a data API?

What are the characteristics of tidy data?

Why are unique observation IDs so important for data cleaning?

# Caveat to Web scraping

I don't expect you to master the tools of web scraping in this course.  
I just want you to know that these things are **possible**, so that you  
**know where to look** in future work.

# Web scraping

Web scraping simply means gathering data from websites.

Last class we learned a particular form of web scraping:  
downloading explicitly structured data files/data APIs.

You can also download information that is not as well structured for  
statistical analysis:

- ▶ HTML tables
- ▶ Text on websites
- ▶ Information that requires you to navigate through web forms

To really master web scraping you need a good knowledge of HTML.



# Key tools

The most basic tools for web scraping in R:

- ▶ `httr`: gather data + simple parsing
- ▶ XML: more advanced parsing
  - ▶ Parsing: the analysis of HTML (and other) markup so that each element is syntactically related in a **parse tree**.

Also take a look at `rvest`. It is a new package that aims to implement features from Python's popular Beautiful Soup.

## Key steps:

1. **Look at** the HTML for the webpage you want to scrape (use Inspect Element in Chrome).
2. **Request** a URL with GET.
3. **Extract** the content from the request with content.
  - ▶ You can extract either the raw text with `as = 'text'` or parse the content with `as = 'parsed'`.
4. **Clean** content (there are many tools for this suited to a variety of problems).

# Web scraping example

Scrape BBC's MP's Expenses table.

HTML markup marks tables using `<table>` tags.

We can use these to extract tabular information and convert them into data frames.

In particular, we want the table tag with the id `expenses_table`.

# Viewing the web pages source

The screenshot displays a web browser window with a table of expenses. The table has columns for names, locations, and various numerical values. A yellow tooltip indicates the table's dimensions as 800px x 19079px. Below the table, the browser's developer tools are open, showing the HTML source code for the table and the CSS styles applied to it.

							Central	Stations	IT prov.	Staff co	Commun.	Travel	TOTAL
Abbott, Ms Diane	LAB	Hackney North & Stoke Newington	0	2,812	20,178	90,325	1,521	3,906	1,351	1,905	7,948	1,789	131,735
Adams, Mr Gerry	SF	West Belfast	21,131	0	21,273	90,278	85	218	1,329	0	0	3,629	137,943

Developer Tools - Elements Panel:

```
<!--endif -->
<table id="expenses_table" class="searchable sortable-enabled">
  <colgroup>_</colgroup>
  <thead>_</thead>
  <tbody>_</tbody>
</table>
<p>* Minus figure denotes repayment as a result of overpayment from previous year</p>
<p class="thdb-no-results-msg">Your search returned no results.</p>
</div>
<script type="text/javascript">_</script>
<!-- E IINC -->
<!-- E B0 -->
<br>
<br>
<br>
<div id="hhrcrm_adsense_middle" class="hhrcrm_adsense hhrcrm_display none">_</div>
```

Developer Tools - Styles Panel:

```
element.style {
}

#expenses_table_wrap
img.headers {
  margin-left: 2px;
  margin-bottom: -2px;
}

media="all"
img {
  border: 1px;
}

media="all"
img, abbr, acronym, fieldset {
```

## Web scraping example

```
library(httr)
library(dplyr)
library(XML)

URL <- 'http://news.bbc.co.uk/2/hi/uk_news/politics/8044207'

# Get and parse all tables on the webpage
tables <- URL %>% GET() %>%
  content(as = 'parsed') %>%
  readHTMLTable()

names(tables)
```

```
## [1] "NULL" "NULL" "NULL"
## [5] "expenses_table" "NULL" "NULL"
```

## Web scraping example

Now we just need to subset the *tables* list for the *expenses\_table* data frame.

```
ExpensesTable <- tables[[5]]
```

```
head(ExpensesTable)[, 1:3]
```

		MP	Party	
## 1	Abbott, Ms Diane	LAB		Hackney North & Stokely
## 2	Adams, Mr Gerry	SF		W
## 3	Afriyie, Adam	CON		
## 4	Ainger, Nick	LAB		Carmarthen West & Pembroke
## 5	Ainsworth, Mr Peter	CON		
## 6	Ainsworth, Rt Hon Bob	LAB		Coventry

# Processing strings

A (frustratingly) large proportion of time web scraping and doing data cleaning generally is taken up with **processing strings**.

**Key tools** for processing strings:

- ▶ knowing your encoding and `iconv` function in base R
- ▶ `grep`, `gsub`, and related functions in base R
- ▶ Regular expressions
- ▶ `stringr` package

# Character encoding: Motivation

Sometimes when you load text into R you will get weird symbols like (the replacement character) or other strange things will happen to the text.

NOTE: remember to always check your data when you import it!

This often happens when R is using the **wrong character encoding**.



# Character encoding

All characters in a computer are **encoded** using some standardised system.

R can recognise latin1 and UTF-8.

- ▶ latin1 is fairly limited (mostly to the latin alphabet)
- ▶ UTF-8 covers a much wider range of characters in many languages

You may need to use the `iconv` function to convert a text to UTF-8 before trying to process it.

# grep, gsub, and related functions

R (and many programming languages) have functions for **identifying** and **manipulating** strings.

# Matching

You can use `grep` and `grepl` to find patterns in a vector.

```
pets <- c('cats', 'dogs', 'a big snake')
```

```
grep(pattern = 'cat', x = pets)
```

```
## [1] 1
```

```
grepl(pattern = 'cat', pets)
```

```
## [1] TRUE FALSE FALSE
```

```
# Subset vector
```

```
pets[grep('cats', pets)]
```

```
## [1] "cats"
```

# Terminology

grep stands for: **G**lobally search a **R**egular **E**xpression and **P**rint

# Manipulation

Use `gsub` to substitute strings.

```
gsub(pattern = 'big', replacement = 'small', x = pets)
```

```
## [1] "cats"          "dogs"          "a small snake"
```

# Regular expressions

Regular expressions are a powerful tool for finding and manipulating strings.

They are special characters that can be used to search for text.

For example:

- ▶ find characters at only the beginning or end of a string
- ▶ find characters that follow or are preceded by a particular character
- ▶ find only the first or last occurrence of a character in a string

Many more possibilities.

# Regular expressions examples

Examples modified from Robin Lovelace.

```
base <- c("cat16_24", "25_34cat", "35_44catch",  
          "45_54Cat", "55_4fat$", 'colour', 'color')
```

```
## Find only all 'cat' regardless of case
```

```
grep('cat', base, ignore.case = T)
```

```
## [1] 1 2 3 4
```

## Regular expressions examples

```
# Find only 'cat' at the end of the string with $  
grep('cat$', base)
```

```
## [1] 2
```

```
# Find only 'cat' at the begining of the string with ^  
grep('^cat', base)
```

```
## [1] 1
```



## Regular expressions examples

```
# Find zero or one of the preceeding character with ?  
grep('colou?r', base)
```

```
## [1] 6 7
```

```
# Find one or more of the preceeding character with +  
grep('colou+r', base)
```

```
## [1] 6
```

```
# Find '$' with the escape character \  
grep('\\$', base)
```

```
## [1] 5
```

## Regular expressions examples

```
# Find string with any single character between 'c' and 'l'  
grep('c.l', base)
```

```
## [1] 6 7
```

```
# Find a range of numbers with [ - ]  
grep('[1-3]', base)
```

```
## [1] 1 2 3
```

```
# Find capital letters  
grep('[A-Z]', base)
```

```
## [1] 4
```

# Simple regular expressions cheatsheet

Character	Use
\$	characters at the end of the string
^	characters at the beginning of the string
?	zero or one of the preceding character
*	zero or more of the preceding character
+	one or more of the preceding character
\	escape character use to find strings that are expressions
.	any single character
[ - ]	a range of characters

# Simple regular expressions cheatsheet

You can also find the cheat-sheet at:  
[SyllabusAndLectures/Lecture7/README](#)

# String processing with stringr

The stringr package has many helpful functions that make dealing with strings a bit **easier**.

## stringr examples

Remove leading and trailing **whitespace** (this can be a real problem when creating consistent variable values):

```
library(stringr)  
  
str_trim(' hello  ')
```

```
## [1] "hello"
```

## stringr examples

**Split** strings (really useful for turning 1 variable into 2):

```
trees <- c('Jomon Sugi', 'Huon Pine')  
  
str_split_fixed(trees, pattern = ' ', n = 2)
```

```
##      [,1]      [,2]  
## [1,] "Jomon"  "Sugi"  
## [2,] "Huon"   "Pine"
```

## More data transformations with dplyr

The **dplyr** package has powerful capabilities to manipulate data frames quickly (many of the functions are written in the compiled language C++).

It is also useful for transforming data from **grouped observations**, e.g. countries, households.



# dplyr

Set up for examples

```
# Create fake grouped data
library(randomNames)
library(dplyr)
library(tidyr)

people <- randomNames(n = 1000)
people <- sort(rep(people, 4))
year <- rep(2010:2013, 1000)
trend_income <- c(30000, 31000, 32000, 33000)
income <- replicate(trend_income + rnorm(4, sd = 20000),
                    n = 1000) %>%
  data.frame() %>%
  gather(obs, value, X1:X1000)
income$value[income$value < 0] <- 0
data <- data.frame(people, year, income = income$value)
```

```
head(data)
```

```
##           people year  income
## 1 Abachiche, Yasmine 2010 39775.83
## 2 Abachiche, Yasmine 2011 32018.14
## 3 Abachiche, Yasmine 2012 43596.81
## 4 Abachiche, Yasmine 2013 50081.46
## 5 Abeyta, Alexandria 2010 14823.64
## 6 Abeyta, Alexandria 2011 31378.52
```

# Simple dplyr

Select rows

```
higher_income <- filter(data, income > 60000)  
  
head(higher_income)
```

##	people	year	income
## 1	Acosta, Richard	2011	60879.85
## 2	Aggrey, Mackenzie James	2011	87091.67
## 3	Ahrens, Jessica	2012	65206.52
## 4	Ajibade, Peter	2013	70418.71
## 5	Alaniz-Soto, Brianna	2012	63102.12
## 6	Almeida Montes, Magnolia	2013	60959.87

# Simple dplyr

Select columns

```
people_income <- select(data, people, income)
```

*# OR*

```
people_income <- select(data, -year)
```

```
head(people_income)
```

```
##           people  income
## 1 Abachiche, Yasmine 39775.83
## 2 Abachiche, Yasmine 32018.14
## 3 Abachiche, Yasmine 43596.81
## 4 Abachiche, Yasmine 50081.46
## 5 Abeyta, Alexandria 14823.64
## 6 Abeyta, Alexandria 31378.52
```

## dplyr with grouped data

Tell dplyr what the groups are in the data with `group_by`.

```
group_data <- group_by(data, people)
head(group_data)[1:5, ]
```

```
## Source: local data frame [5 x 3]
## Groups: people [2]
##
##           people  year  income
##           (fctr) (int)   (dbl)
## 1 Abachiche, Yasmine 2010 39775.83
## 2 Abachiche, Yasmine 2011 32018.14
## 3 Abachiche, Yasmine 2012 43596.81
## 4 Abachiche, Yasmine 2013 50081.46
## 5 Abeyta, Alexandria 2010 14823.64
```

Note: all of the following functions work on **non-grouped data** as well.

## dplyr with grouped data

Now that we have declared the data as grouped, we can do operations on each group.

For example, we can extract the highest and lowest income years for each person:

```
min_max_income <- summarize(group_data,  
                             min_income = min(income),  
                             max_income = max(income))  
head(min_max_income)[1:3, ]
```

```
## Source: local data frame [3 x 3]
```

```
##
```

```
##           people min_income max_income  
##           (fctr)      (dbl)      (dbl)  
## 1 Abachiche, Yasmine 32018.14  50081.46  
## 2 Abeyta, Alexandria 14823.64  31378.52  
## 3      Abraham, Omar 25181.58  30297.51
```

## dplyr with grouped data

We can sort the data using `arrange`.

```
# Sort highest income for each person in ascending order  
ascending <- arrange(min_max_income, max_income)  
head(ascending)[1:3, ]
```

```
## Source: local data frame [3 x 3]  
##  
##           people min_income max_income  
##           (fctr)      (dbl)      (dbl)  
## 1 Anderson, Caren      0.000    6633.351  
## 2 Wreford, Coral       0.000   13942.812  
## 3 Kerr, Israel      7357.525   18625.408
```

## dplyr with grouped data

Add desc to sort in descending order

```
descending <- arrange(min_max_income, desc(max_income))  
head(descending)[1:3, ]
```

```
## Source: local data frame [3 x 3]
```

```
##
```

```
##           people min_income max_income  
##           (fctr)      (dbl)      (dbl)  
## 1 Saunders, Misael      0.000 116819.41  
## 2      Hyman, Aidan      0.000 105891.88  
## 3        Chen, Jerel 9046.961  93060.03
```



## dplyr with grouped data

summarize creates a new data frame with the summarised data.

We can use mutate to add new columns to the original data frame.

```
data <- mutate(group_data,  
               min_income = min(income),  
               max_income = max(income))  
head(data)[1:3, ]
```

```
## Source: local data frame [3 x 5]
```

```
## Groups: people [1]
```

```
##
```

```
##           people  year  income min_income max_income
```

```
##           (fctr) (int)   (dbl)      (dbl)      (dbl)
```

```
## 1 Abachiche, Yasmine 2010 39775.83   32018.14   50081.4
```

```
## 2 Abachiche, Yasmine 2011 32018.14   32018.14   50081.4
```

```
## 3 Abachiche, Yasmine 2012 43596.81   32018.14   50081.4
```

# Seminar: Web scraping and data transformations

**Scrape** and **clean** the Medal Table from <http://www.bbc.com/sport/winter-olympics/2014/medals/countries>.

- ▶ Also, sort by total medals in **descending order**.

Work on **gathering data and cleaning** for **Assignment 3**.