

# DATA607 Assignment 4: Tidying and Transformation Data

*Albert Gilharry*

*15 March 2018*

## Contents

Intro . . . . .	1
Load Libraries . . . . .	1
Preview HTML File . . . . .	2
Load HTML Data & Create Data Frame . . . . .	2
View HTML Data Frame . . . . .	2
Preview XML File . . . . .	4
Load XML Data & Create Data Frame . . . . .	4
View XML Data Table . . . . .	4
Preview JSON File . . . . .	4
Load JSON Data & Create Data Frame . . . . .	4
View JSON Data Table . . . . .	6

## Intro

For this assignment, we were tasked with creating HTML, XML, and JSON files of 3 or our favourite books on one of our favorite topics. At least one of the books must have more than one author. Each of the different file structures should be loaded into R data frames. This is a primer for further work with these structures in the semester.

## Load Libraries

```
library("tidyverse")
```

```
## -- Attaching packages -----  
## v ggplot2 2.2.1      v purrr  0.2.4  
## v tibble  1.4.2      v dplyr  0.7.4  
## v tidyr   0.8.0      v stringr 1.2.0  
## v readr   1.1.1      v forcats 0.2.0  
  
## -- Conflicts -----  
## x dplyr::filter() masks stats::filter()  
## x dplyr::lag()    masks stats::lag()
```

```
library("rvest")
```

```
## Loading required package: xml2  
##  
## Attaching package: 'rvest'  
  
## The following object is masked from 'package:purrr':  
##
```

```
##      pluck
## The following object is masked from 'package:readr':
##
##      guess_encoding
library("XML")

##
## Attaching package: 'XML'
## The following object is masked from 'package:rvest':
##
##      xml
library("methods")
library("jsonlite")

##
## Attaching package: 'jsonlite'
## The following object is masked from 'package:purrr':
##
##      flatten
library("stringr")
library("DT")
library("RCurl")

## Loading required package: bitops
##
## Attaching package: 'RCurl'
## The following object is masked from 'package:tidyr':
##
##      complete
```

## Preview HTML File

### Load HTML Data & Create Data Frame

```
# load HTML data into data frame
url <- "https://raw.githubusercontent.com/albert-gilharry/DATA607-Assignment-5/master/data/books.html"
htmlBooks <- url %>%
  read_html() %>%
  html_nodes("table") %>%
  html_table()

htmlBooks <- htmlBooks[[1]]
```

## View HTML Data Frame

```
datatable(htmlBooks, options = list(filter = FALSE), filter = "none")
```

```
## PhantomJS not found. You can install it with webshot::install_phantomjs(). If it is installed, please
```

```

1 <HTML>
2   <HEAD>
3     <TITLE>My Favorite Books</TITLE>
4   </HEAD>
5   <BODY>
6     <TABLE>
7       <TR>
8         <TH>Title</TH>
9         <TH>Author(s)</TH>
10        <TH>Publisher</TH>
11        <TH>Year</TH>
12        <TH>Pages</TH>
13        <TH>ISBN</TH>
14      </TR>
15      <TR>
16        <TD>The Data Warehouse Toolkit, 3rd Edition</TD>
17        <TD>Ralph Kimball, Margy Ross</TD>
18        <TD>John Wiley & Sons, Inc.</TD>
19        <TD>2013</TD>
20        <TD>563</TD>
21        <TD>978-1-118-53080-1</TD>
22      </TR>
23      <TR>
24        <TD>Data Science of Business, 3rd Edition</TD>
25        <TD>Foster Provost, Tom Fawcett</TD>
26        <TD>O'Reilly Media, Inc.</TD>
27        <TD>2013</TD>
28        <TD>384</TD>
29        <TD>978-1-449-36132-7</TD>
30      </TR>
31      <TR>
32        <TD>Introduction to the Design and Analysis of Algorithms</TD>
33        <TD>Anany Levitin</TD>
34        <TD>Addison-Wesley</TD>
35        <TD>2003</TD>
36        <TD>497</TD>
37        <TD>9780201743951</TD>
38      </TR>
39    </TABLE>
40  </BODY>
41 </HTML>

```

Figure 1: HTML File

## Preview XML File

### Load XML Data & Create Data Frame

The books with multiple authors posed a problem because the built in functionality to convert XML to a data frame concatenates the author nodes without a delimiter. For this reason, I looped through the data to format the authors' names properly. This may not be most efficient way of doing so, but this is a very small data set, so it is fine.

```
url <- getURL("https://raw.githubusercontent.com/albert-gilharry/DATA607-Assignment-5/master/data/books.xml")
doc <- xmlParse(url)

data <- xpathSApply(doc, "//BOOKS/BOOK/AUTHORS",xmlChildren, simplify = TRUE)
authors = c()
for(i in 1:length(data)){
  c <- c()
  for(j in 1:length(data[[i]])){
    c <- append( c, xmlValue(data[[i]][[j]]) )
  }
  authors <- append(authors, paste(unlist(c),collapse = ", "))
}

# use the built in function to create the data frame
xmlBooks <- xmlToDataFrame(url,stringsAsFactors = FALSE)

# fix the authors
xmlBooks$AUTHORS <- authors
```

### View XML Data Table

```
datatable(xmlBooks, options = list(filter = FALSE),filter="none")
```

## Preview JSON File

### Load JSON Data & Create Data Frame

The books with multiple authors posed a problem again because the built in functionality to create a data frame from JSON data attaches an list for the authors. I looped through the data to format the authors. This may not be the most efficient way of doing so, but this is a very small data set, so it is fine.

```
# load JSON data into data frame
url <- getURL("https://raw.githubusercontent.com/albert-gilharry/DATA607-Assignment-5/master/data/books.json")
jsonBooks <- fromJSON(url)
authors = c()
jsonBooks <- jsonBooks$books

# create a comma separated list for authors of each book
for(i in 1:nrow(jsonBooks)){
  authors <- append(authors, paste(unlist( jsonBooks$author[i] ),collapse = ", "))
}
```

```

1  <?xml version="1.0" encoding="UTF-8"?>
2  <BOOKS>
3    <BOOK>
4      <TITLE>The Data Warehouse Toolkit, 3rd Edition</TITLE>
5      <AUTHORS>
6        <AUTHOR>
7          <NAME>Ralph Kimball</NAME>
8        </AUTHOR>
9        <AUTHOR>
10       <NAME>Margy Ross</NAME>
11     </AUTHOR>
12   </AUTHORS>
13   <PUBLISHER>John Wiley & Sons, Inc.</PUBLISHER>
14   <YEAR>2013</YEAR>
15   <PAGES>563</PAGES>
16   <ISBN>978-1-118-53080-1</ISBN>
17 </BOOK>
18 <BOOK>
19   <TITLE>Data Science of Business, 3rd Edition</TITLE>
20   <AUTHORS>
21     <AUTHOR>
22       <NAME>Foster Provost</NAME>
23     </AUTHOR>
24     <AUTHOR>
25       <NAME>Tom Fawcett</NAME>
26     </AUTHOR>
27   </AUTHORS>
28   <PUBLISHER>O'Reilly Media, Inc.</PUBLISHER>
29   <YEAR>2013</YEAR>
30   <PAGES>384</PAGES>
31   <ISBN>978-1-449-36132-7</ISBN>
32 </BOOK>
33 <BOOK>
34   <TITLE>Introduction to the Design and Analysis of Algorithms</TITLE>
35   <AUTHORS>
36     <AUTHOR>
37       <NAME>Anany Levitin</NAME>
38     </AUTHOR>
39   </AUTHORS>
40   <PUBLISHER>Addison-Wesley</PUBLISHER>
41   <YEAR>2003</YEAR>
42   <PAGES>497</PAGES>
43   <ISBN>9780201743951</ISBN>
44 </BOOK>
45 </BOOKS>

```

Figure 2: HTML File

```

1  {
2  "books": [
3      {
4          "title": "The Data Warehouse Toolkit, 3rd Edition",
5          "author": ["Ralph Kimball", "Margy Ross"],
6          "publisher": "John Wiley & Sons, Inc.",
7          "year": 2013,
8          "pages": 563,
9          "isbn": "978-1-118-53080-1"
10     },
11     {
12         "title": "Data Science of Business, 3rd Edition",
13         "author": ["Foster Provost", "Tom Fawcett"],
14         "publisher": "O'Reilly Media, Inc.",
15         "year": 2013,
16         "pages": 384,
17         "isbn": "978-1-449-36132-7"
18     },
19     {
20         "title": "Introduction to the Design and Analysis of Algorithms",
21         "author": ["Anany Levitin"],
22         "publisher": "Addison-Wesley",
23         "year": 2003,
24         "pages": 497,
25         "isbn": "9780201743951"
26     }
27 ]
}

```

Figure 3: HTML File

```

# update authors
jsonBooks$author <- authors

```

## View JSON Data Table

```

datatable(jsonBooks, options = list(filter = FALSE), filter="none")

```