

Use R!

Deborah Nolan  
Duncan Temple Lang

# XML and Web Technologies for Data Sciences with R

 Springer

UseR!

Deborah Nolan  
Duncan Temple Lang

# XML and Web Technologies for Data Sciences with R

# Use R!

*Series Editors:*

Robert Gentleman   Kurt Hornik   Giovanni Parmigiani

For further volumes:

<http://www.springer.com/series/6991>





Deborah Nolan • Duncan Temple Lang

# XML and Web Technologies for Data Sciences with R

 Springer

Deborah Nolan  
Department of Statistics  
University of California  
Berkeley, CA, USA

Duncan Temple Lang  
Department of Statistics  
University of California  
Davis, CA, USA

ISSN 2197-5736 ISSN 2197-5744 (electronic)  
ISBN 978-1-4614-7899-7 ISBN 978-1-4614-7900-0 (eBook)  
DOI 10.1007/978-1-4614-7900-0  
Springer New York Heidelberg Dordrecht London

Library of Congress Control Number: 2013954669

© Springer Science+Business Media New York 2014

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media ([www.springer.com](http://www.springer.com))

*To Doris and Harriet,  
and my teacher Winifred Asprey.*

*— Deborah*

*To Zoë and Suzana,  
and my family farther away.*

*— Duncan*



# Preface

There has been a major change over the last decade in many aspects related to working with data and doing scientific work. As the bloggers at [simplystatistics.org](http://simplystatistics.org) put it, this is a “new era where data are abundant and statisticians are scientists.” The growth of the Web and the numerous data technologies that it has fostered have changed what we can do, and how we do it. It has helped to broaden the focus of statistics from mostly the modeling stage to all stages of data science: finding relevant data, accessing data, reading and transforming data, visualizing the data in rich ways, modeling, and presenting the results and conclusions with compelling, interactive displays. In this book, we describe many technologies and approaches related to all of these additional stages in the data scientist’s work flow. We focus on important and fundamental technologies that are likely to stand the test of time, explain their roles, and compare them with other technologies and approaches. Importantly, we also illustrate how to work with them within the *R* programming environment through numerous comprehensive examples and case studies.

Our focus is on technologies related to the different stages of the data analysis work flow. We can now access so much data from many different sources, in different formats and via many different techniques. We can use formal Web services and application programming interfaces (APIs) or simply scrape data from human-readable Web pages. The data may come in some dialect of *XML*, *HTML* or as a *JSON* document or some other self-describing format. We will explore how we both make Web requests—both simple and sophisticated—and transform the content into data in *R*.

While we can use *R*’s rich graphical functionality, we can also visualize data in new ways with a collection of interactive plots and text in a Web browser, or use applications such as Google Earth to display spatio-temporal data and models. For these, we can use *R* to create these “external” displays as *JavaScript*, *SVG*, or *KML* documents. We can export data from *R* as *JSON* for use in *JavaScript* code to connect the data analysis with the visualization of the results.

Technologies such as Web services and requests, *XML*, and *JSON* are widely used in contexts other than the Web. All of the desktop office suites that provide spreadsheets, word processors, and presentation applications use *XML* to represent the contents of those documents. We also interact with online office suites such as GoogleDocs via authenticated Web requests using *OAuth* and then digest the *XML* content, and we also communicate with newer *NoSQL* databases and other applications using Web service technologies locally. Therefore, the technologies we discuss in this book are, in many ways, fundamental infrastructure for modern computing with data.

One of the important concepts motivating this book is that data analysts and modern statisticians increasingly are working on multi-disciplinary projects or posing their own questions. They need to access data and find auxiliary data, and wrangle them into a usable form. They expect to create rich and informative graphical displays in Web browsers, dashboards, or dynamic reports, and should

be familiar with how to do this. They should also be able to understand the core technologies and know when and how to leverage them. We like to think that we foresaw this evolution when we first started developing the *R* interfaces for these technologies, back in December of 1999—a different century! We are very glad to see the rise of data science as the broadening of statistics and also see the technologies we cover in this book growing in importance in the practice of data analysis. This book is aimed at this new breed of scientist who thinks of statistics as spanning all the stages of the data work flow, and who wants to be involved in all of them.

The book describes a mixture of general infrastructure tools for *R* programmers, e.g., parse an *XML* document or make a Web request, and end-user tools *R* users can employ directly to perform a high-level task, e.g., read an *HTML* table as an *R* data frame or access the contents of a document from Google Docs. The aim in each chapter is to introduce to the reader an important technology and illustrate how it is relevant for modern statisticians. We introduce each of these technologies, provide a succinct overview of the essential concepts and elements, and describe *R*-based tools for working with the particular technology, including numerous packages developed as part of the Omegahat project ([www.omegahat.org](http://www.omegahat.org)). We illustrate each technology with high-level, reasonably simple examples. We also provide more in-depth examples that put the different pieces in context and compare possible approaches. Combining these technologies into a book allows us to incrementally build on the descriptions of the fundamentals.

We have organized this book into three parts. The first part introduces the reader to *XML* and *JSON*. It assumes no prior knowledge of these data formats. It includes discussions of the tools in *R* for reading *XML* and *JSON* files and extracting data from them. It describes various different approaches and computational models for reading *XML* and explains why they are all useful in different circumstances, dealing with documents ranging from small to enormous in size. While we may spend most of our time processing content that was generated by other researchers and Web sites, etc., we often want to create our own *XML* and *JSON* documents, e.g., for displaying plots in Web browsers or Google Earth or upload to Google Docs. We also address this topic in the first part of the book.

Our focus in the second part is on how to obtain data over the Web from remote sites, services, and APIs. These methods include accessing data directly from *HTML* on static Web pages and also from dynamic *HTML* documents via *HTML* forms. We also explore using more structured access via Web services such as *SOAP* (Simple Object Access Protocol), *REST* (Representational State Transfer), and *XML-RPC* (*XML* Remote Procedure Call). Additionally, we may need to obtain data locally from an application running as a server, such as a *NoSQL* database. We will introduce and discuss the common technologies used for these methods. We will also explore the common approaches to authorization and authentication when accessing data from a remote site. Some of these are simple, e.g., passwords sent in a request, and others are more industrial-strength methods using technologies such as *OAuth*.

Accessing data from a remote host involves the client (*R*) communicating with that host. *R* already has several functions to do this. These are sufficient for many common tasks, but as we try to use *R* for more sophisticated network communications and interactions, we need richer tools and a richer, more flexible interface to use those tools. We will discuss the *RCurl* package that gives us these more general facilities. This infrastructure allows us now, and in the future, to harness new technologies built on *HTTP* and other protocols.

The final part of this book presents four in-depth examples that cover: 1) *XML* Schema, a grammar for describing the rules of an *XML* grammar; 2) *SpreadsheetML*, an Office Open *XML* vocabulary for spreadsheets and report writing; 3) Scalable Vector Graphics (*SVG*) for creating interactive graphical displays; and 4) Keyhole Markup Language (*KML*) for displaying geographic data on Google Earth and Maps. In each chapter, we describe the underlying computational model we have developed to work with these *XML* formats to, e.g., programmatically define *R* data structures, create reports, and

design new graphics formats—several important aspects of data science. We hope that these examples serve as case studies for how someone might create an *R* interface to other *XML* vocabularies.

Throughout the book, we present several dozen *R* packages for data scientists to work with Web-related data technologies. Some of these packages are essentially complete and can be useful immediately to *R* users. Others are more infrastructural and/or speculative to illustrate what might be useful and how it might be done. We hope these packages will excite readers to think about how these technologies might be used to do new things and that some readers will want to extend the software to realize new possibilities. For those readers, we have provided ideas for enhancements at the end of many of the chapters and we welcome your contributions and suggestions. If you are interested in extending our work and developing new applications, we encourage you to read about writing *R* packages [145] and useful, slightly advanced aspects of the *R* language in texts such as [30, 56, 108].

We make no claim that *R* is the best tool for working with *XML*, *JSON*, and other Web technologies, or that people should use it in preference to other languages such as *Python*, *Java*, *PERL*, and *JavaScript*. However, we are suggesting that statisticians already using *R* should be able to work with these new technologies without an entire shift in their programming language and environment.

Code from the examples in this book and additional materials are available via the Web page for the book – <http://rxmlwebtech.org>.

## Typographic Conventions

In this book, we discuss numerous different programming languages, and also how we use the languages in combination with one another. We illustrate code in these different languages both inline in the text and as separate code blocks. While the context should clearly indicate the language used, we use color to distinguish the languages and different fonts to distinguish the elements of a language, e.g., to differentiate a function name from a package name in *R*.

The following is a collection of example formats for these languages.

### *R*

A block of *R* code appears as

```
doc = as(filename, "XMLInternalDocument")
xpathSApply(doc, "//section/title/", xmlValue)
```

Output from *R* commands appears as

```
      [,1] [,2] [,3]
[1,]     1     6    11
[2,]     2     7    12
```

And, *R* expressions appear inline as `getNodeSet(x, "//js:*", "js")`.

References to *R* function names appear in the form `xmlParse()`, and names of function parameters look like *filename*. Names of *R* variables appear as *variable*, names of *S4* classes in *R* appear as *XSLStyleSheet*, and *S3* class names as *XMLInternalDocument*. Named elements in an *R* list appear as *name* and slots in an *S4* object as *slot*. Options that we can query and set in *R* via the `options()` function appear as *warning.length*. Also, special values in *R* appear as `TRUE`, `FALSE`, `NULL`, `NA`, `NaN`, etc. Formulas are shown as  $y \sim a + b \mid \text{time}$ . Keywords in the language appear as *while*. We display regular *R* package names as *lattice*, Omegahat packages (i.e., packages distributed from <http://www.omegahat.org>) as *RCurl*, and *Bioconductor* packages as, e.g., *graph*.

## File Names and Directories

We render file names as `filename`, file extensions as **xlsx**, and directories with a trailing `/` to differentiate them from file names, e.g., `/home/frank/`.

## XML

The content for any of the various *XML* vocabularies (e.g., *XHTML*, *KML*, *SVG*, *SpreadsheetML*) is displayed as

```
<?xml version="1.0" encoding="UTF-8"?>
<snapshot>
  <header>
    <total>576803</total>
    <page>1</page>
    <date>2010-01-29T20:00:23Z</date>
    <page_size>1000</page_size>
  </header>
  ...
</snapshot>
```

Inline *XML* content is displayed as `<xs:element name="ARIMA" />`. Element names are shown as `<r:plot>`, while namespace prefixes are rendered as `bioc`. Attributes for an *XML* node are displayed as `href`. *XML* entities appear as `&lt;` or `&amp;`.

## XPath

The code blocks for *XPath* or non-inlined expressions are displayed as

```
/Envelope/Cube/Cube
```

We see an expression inline as `//a[@href]` and names of *XPath* functions appear as `starts-with()`. An *XPath* axis is rendered as **ancestor**, and we show node test expressions as **comment**. We display the literal logical values in *XPath* as `true`.

## JSON

Generic *JSON* content appears as

```
{
  "values": [ 1, 2, 3, 4, 5 ],
  "names": [ "ABC", "DEF", "GHI", "JKL", "MNO" ]
}
```

Inline *JSON* content is shown as arrays with `[1, 2, null, 4]`. Fields in *JSON* content appear as `copyright`. The *JSON* literal logical values are displayed as `true` and `false`. Similarly, the null value appears as `null`. *JSON* number and string values are shown as `3.1415` and `string`.

## JavaScript

a *JavaScript* code block appears as

```
function highlightEdges(evt, row, color)
{
  var labels = edgeTable[row];
  var reset = false;

  if(typeof color == 'undefined') {
    color = "black";
```



```

    reset = true;
  }
}

```

References to *JavaScript* function names appear in the form *getElementById()*, and names of variables appear as *neighbors*. Fields of an object are shown as *myVar* while methods of an object appear as *getValue*. We see inline *JavaScript* expressions as *x = slider.getValue()*.

### RCurl and libcurl Options

The names of curl options appear as, for example, *verbose*.

### HTTP

We refer to different aspects of *HTTP* in various chapters. We represent *HTTP* operations, e.g., **GET** and **DELETE**. We also refer to elements/fields of the *HTTP* header with *ContentType*. We show some or all of the header information from a request or response with

```
GET /folder/docName HTTP/1.1
```

### Shell

A block of shell (sh, bash, csh, tcsh, etc.) code appears as

```
xmllint myDoc.xml
```

Output from shell commands is displayed as

```
149 book.xml
426 springerLatex.xsl
575 total
```

When we refer to a shell command/executable, we see it as *xmllint*. Shell variables are displayed as *XML\_CATALOG\_FILES*. Options or flags for shell commands are rendered as *-noout*.

## Acknowledgements

Many have contributed to this project by reading chapters, suggesting interesting ideas to pursue, and expressing support for our work. We thank them all, including those who participated in our NSF-sponsored Computing in Statistics Workshops. (The material in this book is based upon work supported by the National Science Foundation under Grant No. 0618865.) We especially mention: Gabe Becker, Carl Boettiger, Vince Carey, John Chambers, Robert Gentleman, Jeff Gentry, Tammy Greasby, Kurt Hornik, Ross Ihaka, Cari Kaufman, Bitao Liu, Paul Murrel, Balasubramanian Narasimhan, Karthik Ram, Steve Sein, and Phil Spector. We are also grateful to Daniel Veillard for his *XML* parser and *XSL* toolkit, *libxml2* and *liblibxslt*, and Jonathan Wallace for his *JSON* parser.

We wish to acknowledge the users of our software who have provided useful examples, asked questions that turned into examples, and submitted bug reports. Their contributions have served as inspiration for many of the examples in this book. Open source software improves because of the community of contributors. In that vein, we express our appreciation to the *R* Core for maintaining a valuable infrastructure in which we can explore these technologies.

We thank the students who have taken our computing classes over the past eight years (STAT 133 at Berkeley and STAT 141 and 242 at Davis). They have given us valuable feedback as they explored ideas and used our software when working on their projects and assignments. We also thank the

computing support staff in our departments: Larry Tai at Davis; and Rick Kawin, Ryan Lovett, and the rest of the Statistical Computing Facility staff at Berkeley.

We are grateful to our original editor John Kimmel for his continual encouragement and the new folks at Springer—Jon Gurstelle, Marc Strauss and Hannah Bracken—who supported this project to completion.

Finally, we give special thanks to our families for their patience and encouragement throughout this project.

# Contents

<b>Preface</b>	<b>vii</b>
<b>Part I Data Formats: XML and JSON</b>	<b>1</b>
<b>1 Getting Started with XML and JSON</b>	<b>5</b>
1.1 Introduction	5
1.2 Reading Data from <i>HTML</i> Tables	5
1.3 Reading Data from <i>XML</i> -formatted Documents	8
1.3.1 Extracting Data from <i>XML</i> Attributes	13
1.4 Reading Data from <i>JSON</i> -formatted Documents	14
1.5 Summary of Functions to Read <i>HTML</i> , <i>XML</i> , and <i>JSON</i> into <i>R</i> Data Frames and Lists	17
1.6 Further Reading	17
References	18
<b>2 An Introduction to XML</b>	<b>19</b>
2.1 Overview	19
2.2 Essentials of <i>XML</i>	23
2.2.1 Syntax Checkers	28
2.3 Examples of <i>XML</i> Grammars	29
2.3.1 A Discussion of <i>XML</i> Features	35
2.4 Hierarchical Structure	36
2.5 Additional <i>XML</i> Elements	39
2.6 <i>XML</i> Namespaces	42
2.7 Describing the Structure of Classes of <i>XML</i> Documents: Schema and <i>DTD</i> s	45
2.7.1 The <i>DTD</i>	45
2.7.2 Schema	46
2.8 History of <i>XML</i>	50
2.9 Further Reading	50
References	50
<b>3 Parsing XML Content</b>	<b>53</b>
3.1 Introduction to Reading <i>XML</i> in <i>R</i>	53
3.2 The Document Object Model ( <i>DOM</i> )	54
3.3 Accessing Nodes in the <i>DOM</i>	56
3.4 Parsing Other <i>XML</i> Element Types	63
3.5 Parsing <i>HTML</i> Documents	66
3.6 Reading <i>XML</i> from Different Input Sources	67