

Reading XML

Jeffrey Leek Johns Hopkins Bloomberg School of Public Health

_ ≥ ×

- · Extensible markup language
- · Frequently used to store structured data
- · Particularly widely used in internet applications
- · Extracting XML is the basis for most web scraping
- · Components
- Markup labels that give the text structure
- Content the actual text of the document

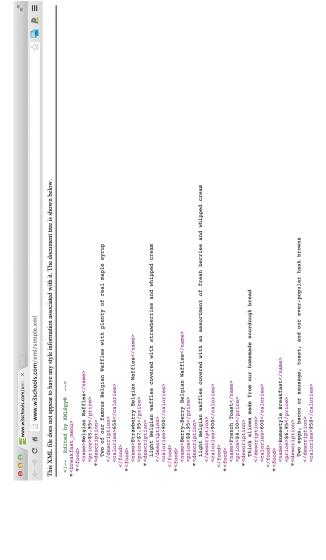
http://en.wikipedia.org/wiki/XML

Tags, elements and attributes

- · Tags correspond to general labels
- Start tags <section>
- End tags </section>
- Empty tags eak />
- · Elements are specific examples of tags
- <Greeting> Hello, world </Greeting>
- Attributes are components of the label
-
- <step number="3"> Connect A to B. </step>

http://en.wikipedia.org/wiki/XML

Example XML file



http://www.w3schools.com/xml/simple.xml

Read the file into R

```
fileUrl <- "http://www.w3schools.com/xml/simple.xml"</pre>
                                                                              doc <- xmlTreeParse(fileUrl,useInternal=TRUE)</pre>
                                                                                                                     rootNode <- xmlRoot(doc)
                                                                                                                                                              xmlName(rootNode)
library(XML)
```

```
[1] "breakfast_menu"
```

names(rootNode)

```
bool bool bool bool bool" "bool" "food" "food" "food" "food" "food"
```

5/14

Directly access parts of the XML document

rootNode[[1]]

```
<description>Two of our famous Belgian Waffles with plenty of real maple syrup/description>
<name>Belgian Waffles</name>
                                                                                                                <calories>650</calories>
                                  </pooj/>
```

rootNode[[1]][[1]]

<name>Belgian Waffles</name>

Programatically extract parts of the file

xmlSApply(rootNode,xmlValue)

Belgian Waffles\$5.95Two of our famous Belgian Waffles with plenty of rea "French Toast\$4.50Thick slices made from our homemade so "Homestyle Breakfast\$6.95Two eggs, bacon or sausage, toast, and our ever-popula "Strawberry Belgian Waffles\$7.95Light Belgian waffles covered with strawberries and "Berry-Berry Belgian Waffles\$8.95Light Belgian waffles covered with an assortment of fresh berries and

7/17

Programatically extract parts of the file

xmlSApply(rootNode,xmlValue)

Belgian Waffles\$5.95Two of our famous Belgian Waffles with plenty of rea "French Toast\$4.50Thick slices made from our homemade sc "Homestyle Breakfast\$6.95Two eggs, bacon or sausage, toast, and our ever-popula "Strawberry Belgian Waffles\$7.95Light Belgian waffles covered with strawberries and "Berry-Berry Belgian Waffles\$8.95Light Belgian waffles covered with an assortment of fresh berries and

8/1/

XPath

- · /node Top level node
- · //node Node at any level
- · node[@attr-name] Node with an attribute name
- · node[@attr-name='bob] Node with attribute name attr-name='bob'

Information from: http://www.stat.berkeley.edu/~statcur/Workshop2/Presentations/XML.pdf

Get the items on the menu and prices

```
xpathSApply(rootNode,"//name",xmlValue)
```

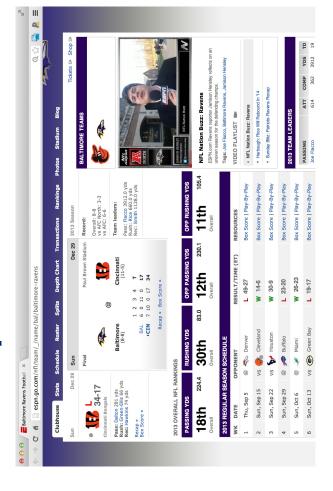
```
[1] "Belgian Waffles"
[4] "French Toast"
```

```
"Strawberry Belgian Waffles" "Berry-Berry Belgian Waffles" "Homestyle Breakfast"
```

```
xpathSApply(rootNode,"//price",xmlValue)
```

```
[1] "$5.95" "$7.95" "$8.95" "$4.50" "$6.95"
```

Another example



http://espn.go.com/nfl/team/_/name/bal/baltimore-ravens

Viewing the source



http://espn.go.com/nfl/team/_/name/bal/baltimore-ravens

```
fileUrl <- "http://espn.go.com/nf1/team/_/name/bal/baltimore-ravens"</pre>
                                                                                                                                                                                                                          teams <- xpathSApply(doc,"//li[@class='team-name']",xmlValue)
                                                                                                                                                 scores <- xpathSApply(doc,"//li[@class='score']",xmlValue)
                                                                            doc <- htmlTreeParse(fileUrl,useInternal=TRUE)</pre>
```

```
"24-18"
              "34-17"
  "19-16"
              "41-7"
              "18-16"
  "19-17"
              "29-26"
  "26-23"
  "23-20"
              "22-20"
              [9] "20-17 OT" "23-20 OT" "19-3"
  "30-9"
  "14-6"
[1] "49–27"
```

teams

```
"Pittsburgh"
 "Green Bay"
               "New York"
"Miami"
                              "New England" "Cincinnati"
              "Chicago"
"Buffalo"
               "Cincinnati"
"Houston"
"Cleveland"
              "Cleveland"
                              "Detroit"
              [7] "Pittsburgh"
                             [13] "Minnesota"
[1] "Denver"
```

13/14

Notes and further resources

- · Official XML tutorials short, long
- · An outstanding guide to the XML package