# Reading XML

Jeffrey Leek

May 18, 2016

# XML

- 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
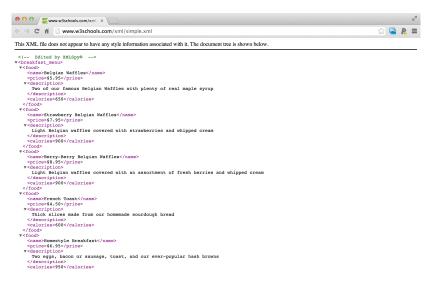- Content - the actual text of the document

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

# Tags, elements and attributes

- Tags correspond to general labels
- Start tags `<section>`
- End tags `</section>`
- Empty tags `<line-break />`
- Elements are specific examples of tags
- `<Greeting>` Hello, world `</Greeting>`
- Attributes are components of the label
- `<img src="jeff.jpg" alt="instructor"/>`
- `<step number="3">` Connect A to B. `</step>`

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

# Example XML file



```
<!-- Edited by XMLSpy® -->
<breakfast_menu>
  <food>
    <name>Belgian Waffles</name>
    <price>$5.95</price>
    <description>
      Two of our famous Belgian Waffles with plenty of real maple syrup
    </description>
    <calories>650</calories>
  </food>
  <food>
    <name>Strawberry Belgian Waffles</name>
    <price>$7.95</price>
    <description>
      Light Belgian waffles covered with strawberries and whipped cream
    </description>
    <calories>900</calories>
  </food>
  <food>
    <name>Berry-Berry Belgian Waffles</name>
    <price>$8.95</price>
    <description>
      Light Belgian waffles covered with an assortment of fresh berries and whipped cream
    </description>
    <calories>900</calories>
  </food>
  <food>
    <name>French Toast</name>
    <price>$4.50</price>
    <description>
      Thick slices made from our homemade sourdough bread
    </description>
    <calories>600</calories>
  </food>
  <food>
    <name>Homestyle Breakfast</name>
    <price>$6.95</price>
    <description>
      Two eggs, bacon or sausage, toast, and our ever-popular hash browns
    </description>
    <calories>950</calories>
  </food>
```

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

# Read the file into R

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

# Directly access parts of the XML document

```
xmlChildren(rootNode)[[1]]

rootNode[[1]]
rootNode[[1]][[1]]
```

# Programatically extract parts of the file
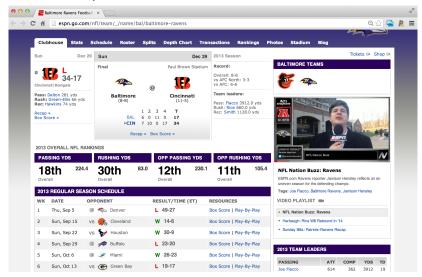
```
xmlSApply(rootNode,xmlValue)
```

# 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)
xpathSApply(rootNode,"//price",xmlValue)
```
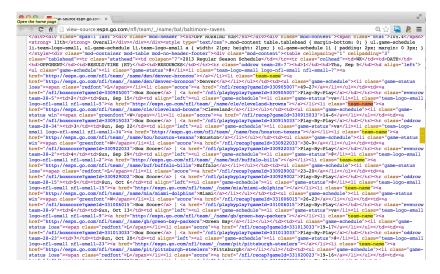
# Another example



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

# Viewing the source



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

# Extract content by attributes

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

# Notes and further resources

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