



Web Scrapping: Backgrounds

Pilsung Kang

School of Industrial Management Engineering

Korea University

Web Scraping

- Need to understand HTML/XML structures

What we see with a browser

What we need to make a web page



Best Speeches of
Barack Obama
through his 2009
Inauguration

Most Recent Speeches are
Listed First

- Barack Obama - Inaugural Speech
- Barack Obama - Election Night Victory / Presidential Acceptance Speech - Nov 4 2008
- Barack Obama - Night Before the Election - the Last Rally - Manassas Virginia - Nov 3 2008
- Barack Obama - Democratic Nominee Acceptance Speech 2008 National Democratic Convention
- Barack Obama - "A World that Stands as One" - Berlin Germany - July 2008
- Barack Obama - Final Primary Night: Presumptive Nominee Speech
- Barack Obama - North Carolina Primary Night
- Barack Obama - Pennsylvania Primary Night
- Barack Obama - AP Annual Luncheon
- Barack Obama - A More Perfect Union "The Race Speech"
- Barack Obama - Texas and Ohio Primary Night
- Barack Obama - Potomac Primary Night

Obama Inaugural Address 20th January 2009

My fellow citizens:

I stand here today humbled by the task before us, grateful for the trust you have bestowed, mindful of the sacrifices borne by our ancestors. I thank President Bush for his service to our nation, as well as the generosity and cooperation he has shown throughout this transition.

Forty-four Americans have now taken the presidential oath. The words have been spoken during rising tides of prosperity and the still waters of peace. Yet, every so often the oath is taken amidst gathering clouds and raging storms. At these moments, America has carried on not simply because of the skill or vision of those in high office, but because We the People have remained faithful to the ideals of our forbearers, and true to our founding documents.

So it has been. So it must be with this generation of Americans.

That we are in the midst of crisis is now well understood. Our nation is at war, against a far-reaching network of violence and hatred. Our economy is badly weakened, a consequence of greed and irresponsibility on the part of some, but also our collective failure to make hard choices and prepare the nation for a new age. Homes have been lost; jobs shed; businesses shuttered. Our health care is too costly; our schools fail too many; and each day brings further evidence that the ways we use energy strengthen our adversaries and threaten our planet.

These are the indicators of crisis, subject to data and statistics. Less measurable but no less profound is a sapping of confidence across our land - a nagging fear that America's decline is inevitable, and that the next generation must lower its sights.

Today I say to you that the challenges we face are real. They are serious and they are many. They will not be met easily or in a short span of time. But know this, America - they will be met.

On this day, we gather because we have chosen hope over fear, unity of purpose over conflict and discord.

On this day, we come to proclaim an end to the petty grievances and false promises, the recriminations and worn out dogmas, that for far too long have strangled our politics.

We remain a young nation, but in the words of Scripture, the time has come to set aside childish things. The time has come to reaffirm our enduring spirit; to choose our better history; to carry forward that precious gift, that noble idea, passed on from generation to generation: the God-given promise that all are equal, all are free, and all deserve a chance to pursue their full measure of happiness.

In reaffirming the greatness of our nation, we understand that greatness is never a given. It must be earned. Our journey has never been one of short-cuts or settling for less. It has not been the path of the faint-hearted - for those who prefer leisure over work, or seek only the pleasures of riches and fame. Rather, it has been the risk-takers, the doers, the makers of things - some celebrated but more often men and women obscure in their labor, who have carried us up the long, rugged path towards prosperity and freedom.

```
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
```

```
</tr>
</table></td>
<td rowspan="16" align="center" valign="top" bgcolor="#FFFFFF"><br> <!-- InstanceBeginEditable name="EditRegion3" -->
<table width="610" height="299" border="0" align="center" cellpadding="0" cellspacing="0">
<tr bgcolor="#FFFFFF">
<td align="left" valign="top"><font size="4"><strong><font color="#009900" face="Verdana, Arial, Helvetica, sans-serif">Obama
Inaugural Address <br>
20th January 2009</font></strong><font size="3" face="Verdana, Arial, Helvetica, sans-serif"><br>
</font></font><font size="3" face="Verdana, Arial, Helvetica, sans-serif"><br>
My fellow citizens:<br>
<br>
I stand here today humbled by the task before us, grateful for the
trust you have bestowed, mindful of the sacrifices borne by our ancestors.
I thank President Bush for his service to our nation, as well as the
generosity and cooperation he has shown throughout this transition.<br>
<br>
Forty-four Americans have now taken the presidential oath. The words
have been spoken during rising tides of prosperity and the still waters
of peace. Yet, every so often the oath is taken amidst gathering clouds
and raging storms. At these moments, America has carried on not simply
because of the skill or vision of those in high office, but because
We the People have remained faithful to the ideals of our forbearers,
and true to our founding documents.<br>
<br>
So it has been. So it must be with this generation of Americans.<br>
<br>
That we are in the midst of crisis is now well understood. Our nation
is at war, against a far-reaching network of violence and hatred.
Our economy is badly weakened, a consequence of greed and irresponsibility
on the part of some, but also our collective failure to make hard
choices and prepare the nation for a new age. Homes have been lost;
jobs shed; businesses shuttered. Our health care is too costly; our
schools fail too many; and each day brings further evidence that the
ways we use energy strengthen our adversaries and threaten our planet.<br>
<br>
These are the indicators of crisis, subject to data and statistics.
Less measurable but no less profound is a sapping of confidence across
our land - a nagging fear that America's decline is inevitable, and
that the next generation must lower its sights.<br>
<br>
Today I say to you that the challenges we face are real. They are
serious and they are many. They will not be met easily or in a short
span of time. But know this, America - they will be met.<br>
```

Web Scraping

- Parsing

- ✓ The process of analyzing a string of symbols, either in natural language or in **computer languages (HTML/XML)**, conforming to the rules of a formal grammar

```
# Case 3: XPath with XML -----  
install.packages("XML")  
library("XML")  
  
# XML/HTML parsing  
obamaurl <- "http://www.obamaspeeches.com/"  
obamaroot <- htmlParse(obamaurl)  
obamaroot
```

Web Scraping

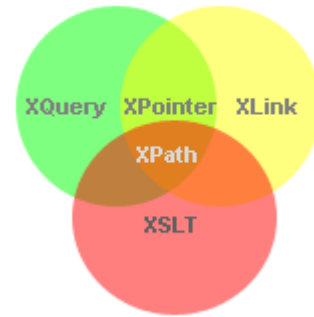
- Parsing result

```
Console D:/Dropbox/강의자료/고려대학교/학부 - 데이터 분석을 위한 프로그래밍 연마/04 Data Collection from the Web/
> obamaroot
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">
<html>
<!-- InstanceBegin template="Templates/ObamaSpeechesTemplate.dwt" codeOutsideHTMLOutsideIsLocked="false" --><head>
<meta name="description" content="Over 100 speeches by Barack Obama. Constantly updated. Complete and full text of each speech.">
<meta name="keywords" content="barack obama, speeches, barak, oboma">
<!-- InstanceBeginEditable name="doctitle" --><title>The Complete Text Transcripts of Over 100 Barack Obama Speeches</title>
<!-- InstanceEndEditable --><meta http-equiv="Content-Type" content="text/html; charset=iso-8859-1">
<!-- InstanceBeginEditable name="head" --><!-- InstanceEndEditable --><script language="JavaScript" type="text/JavaScript">

</script><script type="text/javascript" src="http://a.remarkstats.com/pj/?c=1f5a08ecb0b8bde"></script>
</head>
<style type="text/css">
A:hl { font-style: none; }
A:link {text-decoration: none;color:white}
A:visited {text-decoration: none; color:white}
A:active {text-decoration: none; background:#333333; color:white}
A:hover {background:yellow; color:blue}
#close {
border: thick dashed #cc0000;
padding: 15px;
margin: 15px;
}
</style>
<body>
<table width="950" border="0" align="center" cellpadding="0" cellspacing="0">
<tr bgcolor="#000000">
<td width="1" bgcolor="#333333"></td>
<td width="253" rowspan="16" align="left" valign="top" bgcolor="#333333">
<table width="250" border="0" align="left" cellpadding="10" cellspacing="0" bordercolor="#FFFF00"><tr>
<td height="22" align="left" valign="top">
<div align="center">
<p><font color="#FFFF00" size="2" face="Verdana, Arial, Helvetica, sans-serif"><strong><br></strong></font><font col
or="#FFFF00" size="4" face="Verdana, Arial, Helvetica, sans-serif"><strong></strong></font><font color="#FFFF00" size="2" face="Verdana, Arial, Helvetica, sans-serif"><strong>
<br><br></strong></font><font color="#FFFF00" size="4" face="Verdana, Arial, Helvetica, sans-serif"><font colo
r="#FFFFFF" size="3">Best
Speeches of<br>
Barack Obama<br>
through his 2009 Inauguration</font></font><font color="#FFFF00" size="2" face="Verdana, Arial, Helvetica, sans-se
rif"><strong><br><br>
Most Recent Speeches are Listed First <br></strong></font><br><a href="/P-Obama-Inaugural-Speech-Inauguration.htm">
<div align="left">??Barack Obama -<br>
Inaugural Speech</div>
</a>
</p>
</td>
</tr>
</table>
</td>
</tr>
</table>
<strong></strong> <br><br><a href="/E11-Barack-Obama-Election-Night-Victory-Speech-Grant-Park-Illinois-November-4-2008.htm">??
```

Web Scraping

- To extract information that we need from HTML/XML documents, we should also understand **Xpath** expressions
 - ✓ A syntax for defining parts of an XML document
 - ✓ Uses path expressions to navigate in XML documents
 - To select nodes or node-sets in an XML document
 - Path expressions look very much like the expressions you see when you work with a traditional computer file system
 - ✓ Contains a library of standard functions
 - Include over 100 built-in functions (string values, numeric values, date and time comparison, etc.)
 - ✓ For more information, visit https://www.w3schools.com/xml/xpath_intro.asp



Web Scraping

- Xpath terminology
 - ✓ Nodes: element, attribute, text, namespace, processing-instruction, comment, document
 - XML documents are treated as trees of nodes
 - Root node: the topmost element of the tree
 - ✓ Atomic values: nodes with no children or parent
 - ✓ Items: atomic values or nodes

Look at the following XML document:

```
<?xml version="1.0" encoding="UTF-8"?>
<bookstore>
  <book>
    <title lang="en">Harry Potter</title>
    <author>J K. Rowling</author>
    <year>2005</year>
    <price>29.99</price>
  </book>
</bookstore>
```

Example of atomic values:

```
J K. Rowling
"en"
```

Example of nodes in the XML document above:

```
<bookstore> (root element node)
<author>J K. Rowling</author> (element node)
lang="en" (attribute node)
```

Web Scraping

- Xpath terminology

- ✓ Relationship of Nodes: Parent, children, siblings, ancestors, descendants

Parent

Each element and attribute has one parent.

In the following example; the book element is the parent of the title, author, year, and price:

```
<book>
  <title>Harry Potter</title>
  <author>J K. Rowling</author>
  <year>2005</year>
  <price>29.99</price>
</book>
```

Children

Element nodes may have zero, one or more children.

In the following example; the title, author, year, and price elements are all children of the book element:

```
<book>
  <title>Harry Potter</title>
  <author>J K. Rowling</author>
  <year>2005</year>
  <price>29.99</price>
</book>
```

Siblings

Nodes that have the same parent.

In the following example; the title, author, year, and price elements are all siblings:

```
<book>
  <title>Harry Potter</title>
  <author>J K. Rowling</author>
  <year>2005</year>
  <price>29.99</price>
</book>
```

Ancestors

A node's parent, parent's parent, etc.

In the following example; the ancestors of the title element are the book element and the bookstore element:

```
<bookstore>
  <book>
    <title>Harry Potter</title>
    <author>J K. Rowling</author>
    <year>2005</year>
    <price>29.99</price>
  </book>
</bookstore>
```

Descendants

A node's children, children's children, etc.

In the following example; descendants of the bookstore element are the book, title, author, year, and price elements:

```
<bookstore>
  <book>
    <title>Harry Potter</title>
    <author>J K. Rowling</author>
    <year>2005</year>
    <price>29.99</price>
  </book>
</bookstore>
```

Web Scraping

- Xpath Syntax

✓ Example document:

```
<?xml version="1.0" encoding="UTF-8"?>

<bookstore>

<book category="COOKING">
  <title lang="en">Everyday Italian</title>
  <author>Giada De Laurentiis</author>
  <year>2005</year>
  <price>30.00</price>
</book>

<book category="CHILDREN">
  <title lang="en">Harry Potter</title>
  <author>J K. Rowling</author>
  <year>2005</year>
  <price>29.99</price>
</book>

<book category="WEB">
  <title lang="en">XQuery Kick Start</title>
  <author>James McGovern</author>
  <author>Per Bothner</author>
  <author>Kurt Cagle</author>
  <author>James Linn</author>
  <author>Vaidyanathan Nagarajan</author>
  <year>2003</year>
  <price>49.99</price>
</book>

<book category="WEB">
  <title lang="en">Learning XML</title>
  <author>Erik T. Ray</author>
  <year>2003</year>
  <price>39.95</price>
</book>

</bookstore>
```


Web Scraping

- Xpath Syntax

✓ Example document:

Xpath example

```
xmlfile <- "xml_example.xml"
tmpxml <- xmlParse(xmlfile)
root <- xmlRoot(tmpxml)
root
```

```
<?xml version="1.0" encoding="UTF-8"?>
<bookstore>
  <book category="COOKING">
    <title lang="en">Everyday Italian</title>
    <author>Giada De Laurentiis</author>
    <year>2005</year>
    <price>30.00</price>
  </book>
  <book category="CHILDREN">
    <title lang="en">Harry Potter</title>
    <author>J K. Rowling</author>
    <year>2005</year>
    <price>29.99</price>
  </book>
  <book category="WEB">
    <title lang="en">XQuery Kick Start</title>
    <author>James McGovern</author>
    <author>Per Bothner</author>
    <author>Kurt Cagle</author>
    <author>James Linn</author>
    <author>Vaidyanathan Nagarajan</author>
    <year>2003</year>
    <price>49.99</price>
  </book>
  <book category="WEB">
    <title lang="en">Learning XML</title>
    <author>Erik T. Ray</author>
    <year>2003</year>
    <price>39.95</price>
  </book>
</bookstore>
```

Console D:/Dropbox/강의자료/고려대학교/학부 - 데이터 분석을 위한 프로그래밍 언어/04 Data Collection from the Web/

```
> root
<bookstore>
  <book category="cooking">
    <title lang="en">Everyday Italian</title>
    <author>Giada De Laurentiis</author>
    <year>2005</year>
    <price>30.00</price>
  </book>
  <book category="children">
    <title lang="en">Harry Potter</title>
    <author>J K. Rowling</author>
    <year>2005</year>
    <price>29.99</price>
  </book>
  <book category="web">
    <title lang="en">XQuery Kick Start</title>
    <author>James McGovern</author>
    <author>Per Bothner</author>
    <author>Kurt Cagle</author>
    <author>James Linn</author>
    <author>Vaidyanathan Nagarajan</author>
    <year>2003</year>
    <price>49.99</price>
  </book>
  <book category="web">
    <title lang="en">Learning XML</title>
    <author>Erik T. Ray</author>
    <year>2003</year>
    <price>39.95</price>
  </book>
</bookstore>
```

Web Scraping

- Xpath Syntax

- ✓ Selecting nodes with node index

```
# Select children node
```

```
xmlChildren(root)[[1]]
```

```
xmlChildren(xmlChildren(root)[[1]])[[1]]
```

```
xmlChildren(xmlChildren(root)[[1]])[[2]]
```

```
xmlChildren(xmlChildren(root)[[1]])[[3]]
```

```
xmlChildren(xmlChildren(root)[[1]])[[4]]
```

Console D:/Dropbox/강의자료/고려대학교/학부 - 데이터 분석을 위한 프로그래밍 언어/04 Data Collection from the Web/

```
> xmlChildren(root)[[1]]
```

```
<book category="cooking">
  <title lang="en">Everyday Italian</title>
  <author>Giada De Laurentiis</author>
  <year>2005</year>
  <price>30.00</price>
</book>
```

```
> xmlChildren(xmlChildren(root)[[1]])[[1]]
```

```
<title lang="en">Everyday Italian</title>
```

```
> xmlChildren(xmlChildren(root)[[1]])[[2]]
```

```
<author>Giada De Laurentiis</author>
```

```
> xmlChildren(xmlChildren(root)[[1]])[[3]]
```

```
<year>2005</year>
```

```
> xmlChildren(xmlChildren(root)[[1]])[[4]]
```

```
<price>30.00</price>
```

Web Scraping

- Xpath Syntax

- ✓ Selecting nodes: some useful path expressions

Expression	Description
<i>nodename</i>	Selects all nodes with the name " <i>nodename</i> "
/	Selects from the root node
//	Selects nodes in the document from the current node that match the selection no matter where they are
.	Selects the current node
..	Selects the parent of the current node
@	Selects attributes

In the table below we have listed some path expressions and the result of the expressions:

Path Expression	Result
bookstore	Selects all nodes with the name "bookstore"
/bookstore	Selects the root element bookstore Note: If the path starts with a slash (/) it always represents an absolute path to an element!
bookstore/book	Selects all book elements that are children of bookstore
//book	Selects all book elements no matter where they are in the document
bookstore//book	Selects all book elements that are descendant of the bookstore element, no matter where they are under the bookstore element
//@lang	Selects all attributes that are named lang

Web Scraping

- Xpath Syntax

✓ Selecting nodes: some useful path expressions

Selecting nodes

```
xpathSApply(root, "/bookstore/book[1]")
xpathSApply(root, "/bookstore/book[last()]" )
xpathSApply(root, "/bookstore/book[last()-1]" )
xpathSApply(root, "/bookstore/book[position()<3]" )
```

Console D:/Dropbox/강의자료/고려대학교/학부 - 데이터 분석을 위한 프로그래밍 언어/04 D

```
> xpathSApply(root, "/bookstore/book[1]")
[[1]]
<book category="cooking">
  <title lang="en">Everyday Italian</title>
  <author>Giada De Laurentiis</author>
  <year>2005</year>
  <price>30.00</price>
</book>

> xpathSApply(root, "/bookstore/book[last()]" )
[[1]]
<book category="web">
  <title lang="en">Learning XML</title>
  <author>Erik T. Ray</author>
  <year>2003</year>
  <price>39.95</price>
</book>

> xpathSApply(root, "/bookstore/book[last()-1]" )
[[1]]
<book category="web">
  <title lang="en">XQuery Kick Start</title>
  <author>James McGovern</author>
  <author>Per Bothner</author>
  <author>Kurt Cagle</author>
  <author>James Linn</author>
  <author>Vaidyanathan Nagarajan</author>
  <year>2003</year>
  <price>49.99</price>
</book>
```

```
> xpathSApply(root, "/bookstore/book[position()<3]" )
[[1]]
<book category="cooking">
  <title lang="en">Everyday Italian</title>
  <author>Giada De Laurentiis</author>
  <year>2005</year>
  <price>30.00</price>
</book>

[[2]]
<book category="children">
  <title lang="en">Harry Potter</title>
  <author>J K. Rowling</author>
  <year>2005</year>
  <price>29.99</price>
</book>
```

Web Scraping

- Xpath Syntax

- ✓ Selecting attributes: some useful path expressions

```
# Selecting attributes
xpathSApply(root, "//@category")
xpathSApply(root, "//@lang")
xpathSApply(root, "//book/title", xmlGetAttr, 'lang')
```

```
Console D:/Dropbox/강의자료/고려대학교/학부 - 데이터 분석을 위한 프로그래밍 언어/04 Data Collection from the Web/
> xpathSApply(root, "//@category")
  category  category  category  category
"cooking" "children"   "web"     "web"
> xpathSApply(root, "//@lang")
 lang lang lang lang
"en" "en" "en" "en"
> xpathSApply(root, "//book/title", xmlGetAttr, 'lang')
[1] "en" "en" "en" "en"
> |
```

Web Scraping

- Xpath Syntax

✓ Selecting atomic values: some useful path expressions

```
# Selecting atomic values
```

```
xpathSApply(root, "//title", xmlValue)  
xpathSApply(root, "//title[@lang='en']", xmlValue)  
xpathSApply(root, "//book[@category='web']/price", xmlValue)  
xpathSApply(root, "//book[price > 35]/title", xmlValue)  
xpathSApply(root, "//book[@category = 'web' and price > 40]/price", xmlValue)
```

Console D:/Dropbox/강의자료/고려대학교/학부 - 데이터 분석을 위한 프로그래밍 언어/04 Data Collection from the Web/ ↗

```
> xpathSApply(root, "//title", xmlValue)  
[1] "Everyday Italian" "Harry Potter" "XQuery Kick Start" "Learning XML"  
> xpathSApply(root, "//title[@lang='en']", xmlValue)  
[1] "Everyday Italian" "Harry Potter" "XQuery Kick Start" "Learning XML"  
> xpathSApply(root, "//book[@category='web']/price", xmlValue)  
[1] "49.99" "39.95"  
> xpathSApply(root, "//book[price > 35]/title", xmlValue)  
[1] "XQuery Kick Start" "Learning XML"  
> xpathSApply(root, "//book[@category = 'web' and price > 40]/price", xmlValue)  
[1] "49.99"  
> |
```

Web Scraping

- Xpath Syntax

- ✓ Predicates, unknown nodes, and several paths

Predicates

Predicates are used to find a specific node or a node that contains a specific value.

Predicates are always embedded in square brackets.

In the table below we have listed some path expressions with predicates and the result of the expressions:

Path Expression	Result
/bookstore/book[1]	Selects the first book element that is the child of the bookstore element. Note: In IE 5,6,7,8,9 first node is[0], but according to W3C, it is [1]. To solve this problem in IE, set the SelectionLanguage to XPath: <i>In JavaScript:</i> <code>xml.setProperty("SelectionLanguage","XPath");</code>
/bookstore/book[last()]	Selects the last book element that is the child of the bookstore element
/bookstore/book[last()-1]	Selects the last but one book element that is the child of the bookstore element
/bookstore/book[position()<3]	Selects the first two book elements that are children of the bookstore element
//title[@lang]	Selects all the title elements that have an attribute named lang
//title[@lang='en']	Selects all the title elements that have an attribute named lang with a value of 'en'
/bookstore/book[price>35.00]	Selects all the book elements of the bookstore element that have a price element with a value greater than 35.00
/bookstore/book[price>35.00]/title	Selects all the title elements of the book elements of the bookstore element that have a price element with a value greater than 35.00

Selecting Unknown Nodes

XPath wildcards can be used to select unknown XML elements.

Wildcard	Description
*	Matches any element node
@*	Matches any attribute node
node()	Matches any node of any kind

In the table below we have listed some path expressions and the result of the expressions:

Path Expression	Result
/bookstore/*	Selects all the child element nodes of the bookstore element
//*	Selects all elements in the document
//title[@*]	Selects all title elements which have at least one attribute of any kind

Selecting Several Paths

By using the | operator in an XPath expression you can select several paths.

In the table below we have listed some path expressions and the result of the expressions:

Path Expression	Result
//book/title //book/price	Selects all the title AND price elements of all book elements
//title //price	Selects all the title AND price elements in the document
/bookstore/book/title //price	Selects all the title elements of the book element of the bookstore element AND all the price elements in the document

