



Web Data Collection with R

Session 5: Collecting Data from Static Websites: XPath

Sascha Göbel
sascha.goebel@uni-konstanz.de

December 2nd, 2020

Today: XPath

- identify information in HTML/XML
- download target HTML/XML
- parse the HTML/XML file
- extract the actual information using XPath queries \Leftarrow

XPath

XML path language

- query language for extracting parts from XML/HTML files
- using characteristics and relations of nodes
- a “node” is the entire information encapsulated in tags
- a “node set” is a selection of multiple nodes

Why XPath?

- designed for XML-style files
- intuitive
- flexible
- efficient
- robust

XPath

XML path language

- query language for extracting parts from XML/HTML files
- using characteristics and relations of nodes
- a “node” is the entire information encapsulated in tags
- a “node set” is a selection of multiple nodes

Why XPath?

- designed for XML-style files
- intuitive
- flexible
- efficient
- robust

Queries

XPath query fundamentals

- based on hierarchical order of nodes
- expressed by sequence of nodes leading to the target node
- literally depicting a location path using /
- a character string that reads from left to right
- absolute path:
 - start from root node
 - proceed through all nodes on the way to the target
 - e.g., “/html/body/div/p/i”
- relative path:
 - jump between nodes using “//”
 - e.g., “//p/i”
 - computationally more expensive

XPath query fundamentals

- wildcard “*” matches any node
- “.” selects the current node
- “..” selects the node one level up
- “@” selects attributes
- separate multiple queries by “|” or pass as a vector
- queries have no length limit

How to write a query?

- define where the actual information is located in the DOM
- develop expression from there on, up the tree
- usually a trial-and-error process

XPath query fundamentals

- wildcard “*” matches any node
- “.” selects the current node
- “..” selects the node one level up
- “@” selects attributes
- separate multiple queries by “|” or pass as a vector
- queries have no length limit

How to write a query?

- define where the actual information is located in the DOM
- develop expression from there on, up the tree
- usually a trial-and-error process

Using node relations

- we can identify nodes based on their relation to other nodes
- using notation based on family relationships
- describing the relation of node2 to node1
- “node1/relation::node2”

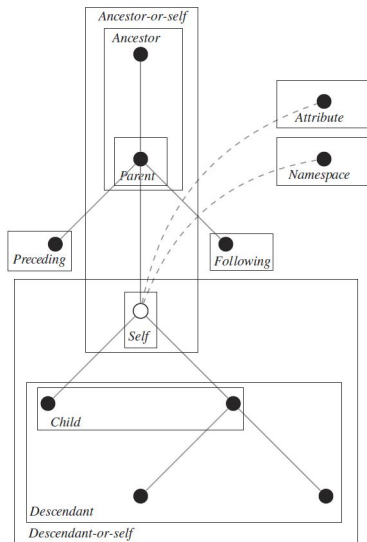


Figure 4.2 Visualizing node relations. Descriptions are presented in relation to the white node

Source: Munzert et al. 2015

Table 4.1 XPath axes

Axis name	Result
ancestor	Selects all ancestors (parent, grandparent, etc.) of the current node
ancestor-or-self	Selects all ancestors (parent, grandparent, etc.) of the current node and the current node itself
attribute	Selects all attributes of the current node
child	Selects all children of the current node
descendant	Selects all descendants (children, grandchildren, etc.) of the current node
descendant-or-self	Selects all descendants (children, grandchildren, etc.) of the current node and the current node itself
following	Selects everything in the document after the closing tag of the current node
following-sibling	Selects all siblings after the current node
namespace	Selects all namespace nodes of the current node
parent	Selects the parent of the current node
preceding	Selects all nodes that appear before the current node in the document except ancestors, attribute nodes, and namespace nodes
preceding-sibling	Selects all siblings before the current node
self	Selects the current node

Source: http://www.w3schools.com/xpath/xpath_axes.asp

Source: Munzert et al. 2015

Using node predicates

- functions applied to node name, value, or attribute
- evaluate whether a condition is true (selected) or false (not selected)
- “node1[predicate]”
- work in conjunction with operators
- when nesting predicates swap “=” with “,”

Table 4.2 Overview of some important XPath functions

Function	Description	Example
<code>name(<node>)</code>	Returns the name of <node> or the first node in a node set	<code>//*[name()='title'];</code> Returns: <title>
<code>text(<node>)</code>	Returns the value of <node> or the first node in a node set	<code>//*[text()='The book homepage'];</code> Returns: <i> with value <i>The book homepage</i>
<code>@attribute</code>	Returns the value of a node's attribute or of the first node in a node set	<code>//div[@id='R_Inventor'];</code> Returns: <div> with attribute <i>id</i> value <i>R_Inventor</i>
<code>string-length(str1)</code>	Returns the length of str1. If there is no string argument, it returns the length of the string value of the current node	<code>//h1[string-length()>11];</code> Returns: <h1> with value <i>Robert Gentleman</i>
<code>translate(str1, str2, str3)</code>	Converts str1 by replacing the characters in str2 with the characters in str3	<code>//div[translate(/@date, '2003', '2005')='June/2005'];</code> Returns: first <div> node with date attribute value <i>June/2003</i>
<code>contains(str1, str2)</code>	Returns TRUE if str1 contains str2, otherwise FALSE	<code>//div[contains(@id, 'Inventor')];</code> Returns: first <div> node with id attribute value <i>R_Inventor</i>
<code>starts-with(str1, str2)</code>	Returns TRUE if str1 starts with str2, otherwise FALSE	<code>//i[starts-with(text(), 'The')];</code> Returns: <i> with value <i>The book homepage</i>
<code>substring-before(str1, str2)</code>	Returns the start of str1 before str2 occurs in it	<code>//div[substring-before(@date, '/')='June'];</code> Returns: <div> with date attribute value <i>June/2003</i>
<code>substring-after(str1, str2)</code>	Returns the remainder of str1 after str2 occurs in it	<code>//div[substring-after(@date, '/')=2003];</code> Returns: <div> with date attribute value <i>June/2003</i>
<code>not(arg)</code>	Returns TRUE if the boolean value is FALSE, and FALSE if the boolean value is TRUE	<code>//div[not(contains(@id, 'Inventor'))];</code> Returns: the <div> node that does not contain the string <i>Inventor</i> in its id attribute value
<code>local-name(<node>)</code>	Returns the name of the current <node> or the first node in a node set—without the namespace prefix	<code>//*[local-name()='address'];</code> Returns: <address>
<code>count(<node>)</code>	Returns the count of a nodeset <node>	<code>//div[count(/@a)=0];</code> Result: The second <div> with one <a> child
<code>position(<node>)</code>	Returns the index position of <node> that is currently being processed	<code>//div/p[position()=1];</code> Result: The first <p> node in each <div> node
<code>last()</code>	Returns the number of items in the processed node list <node>	<code>//div/p[last()];</code> Result: The last <p> node in each <div> node

Source: Munzert et al. 2015

Table 4.3 XPath operators

Operators	Description	Example
	Computes two node sets	//i //b
+	Addition	5 + 3
-	Subtraction	8 - 2
*	Multiplication	8 * 5
div	Division	8 div 5
=	Equal	count = 27
!=	Not equal	count != 27
<	Less than	count < 27
≤	Less than or equal to	count ≤ 27
>	Greater than	count > 27
≥	Greater than or equal to	count ≥ 27
or	Or	count = 27 or count = 28
and	And	count > 26 and count < 30
mod	Modulo (division remainder)	7 mod 2

Source: Adapted from http://www.w3schools.com/xpath/xpath_operators.asp

Source: Munzert et al. 2015

XPath in R

R package “rvest”

- implements XPath 1.0
- no regex support
- workhorse functions:
 - `html_node()` - return first node matching the query
 - `html_nodes()` - return all nodes matching the query
 - `html_text()` - extract the value of a node
 - `html_attr()` - extract the specified attribute of a node
 - `html_table()` - parse a table into a data frame

Helpers

Tools for query construction

- Copy XPath from browser inspection tool
- Selector gadget extension
- Best to learn and code XPath yourself!

Next Session

XPath practice session

- please bring an idea and URL for web data collection
- the website should be static and require XPath, not (yet) API access
 - no Twitter!
- apply and practice what you learned thus far
- chance to solve problems and ask questions in class