

Informatics 1: Data & Analysis

Lecture 11: Navigating XML using XPath

Ian Stark

School of Informatics
The University of Edinburgh

Tuesday 26 February 2013
Semester 2 Week 6



Lecture Timing

This is Inf1-DA Lecture 11, in Week 6.

There is no Inf1-DA lecture on Friday, 1 March.

Inf1-DA Lecture 12 is on Tuesday 5 March, in Week 7.

Normal service then resumes.

XML

We start with technologies for modelling and querying *semistructured data*.

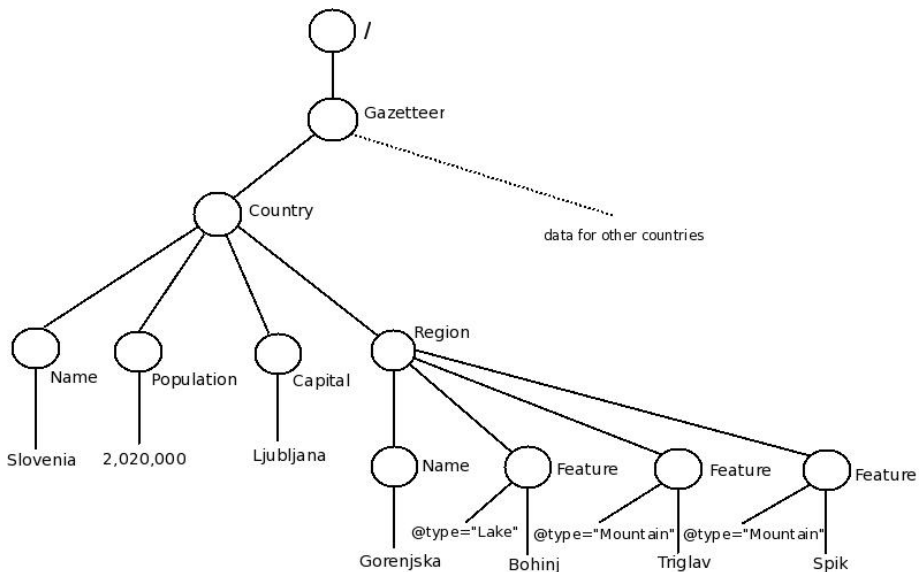
- Semistructured Data: Trees and XML
- Schemas for structuring XML
- Navigating and querying XML with XPath

Corpora

One particular kind of semistructured data is large bodies of written or spoken text: each one a *corpus*, plural *corpora*.

- Corpora: What they are and how to build them
- Applications: corpus analysis and data extraction

Sample Semistructured Data



Sample Semistructured Data in XML

```
<Gazetteer>
  <Country>
    <Name>Slovenia</Name>
    <Population>2,020,000</Population>
    <Capital>Ljubljana</Capital>
    <Region>
      <Name>Gorenjska</Name>
      <Feature type="Lake">Bohinj</Feature>
      <Feature type="Mountain">Triglav</Feature>
      <Feature type="Mountain">Spik</Feature>
    </Region>
  </Country>
  <!-- data for other countries here -->
</Gazetteer>
```

How to Extract Information from an XML Document?

Since an XML document is a text document, we could simply use conventional text search to look for data.

However, this ignores all the document structure.

A more powerful approach is to use a dedicated language for forming queries based on the tree structure of an XML document.

This is (yet another) *domain-specific language*.

With such a language we can, for example:

- Perform database-style queries on data published as XML;
- Extract annotated content from marked-up text documents;
- Identify information captured in the tree structure itself.

XQuery and XPath

XQuery is a powerful declarative query language for extracting information from XML documents.

As well as using XML documents for its source data, XQuery can also produce XML documents as output, so we can view it as an XML *transformation* language.

However, the XQuery language is complex, and we shall not investigate it further here.

XPath is a sublanguage of XQuery, used for navigating XML documents using *path expressions*.

XPath can be viewed as a rudimentary query language in its own right.

It is also an important component of other XML application languages (XML Schema, XSLT, XForms, ...).

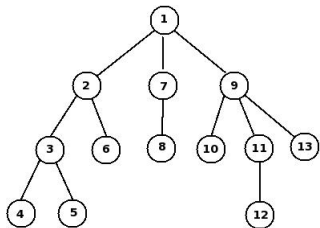
XPath Location Paths

An XPath *location path* (or *path expression*) identifies a set of nodes within an XML document tree.

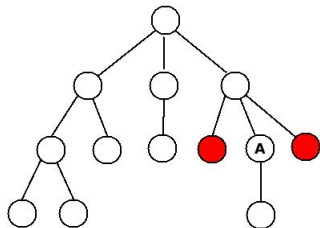
- The location path describes a set of possible paths from the root of the tree.
- The set of nodes identified is all those reached as final destinations of these paths.

When using a location path as a query on a document, this set of nodes is returned as a list (without duplicates) sorted in *document order* — the order the nodes appeared in the original XML document.

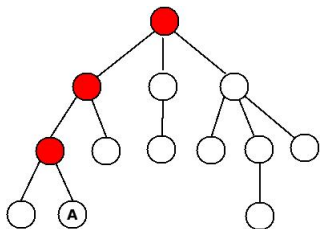
Family Tree Navigation



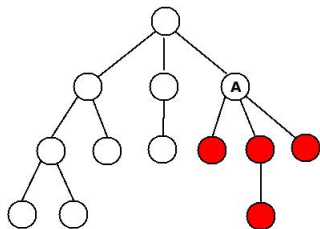
Document order



Siblings of A



Ancestors of A



Descendants of A

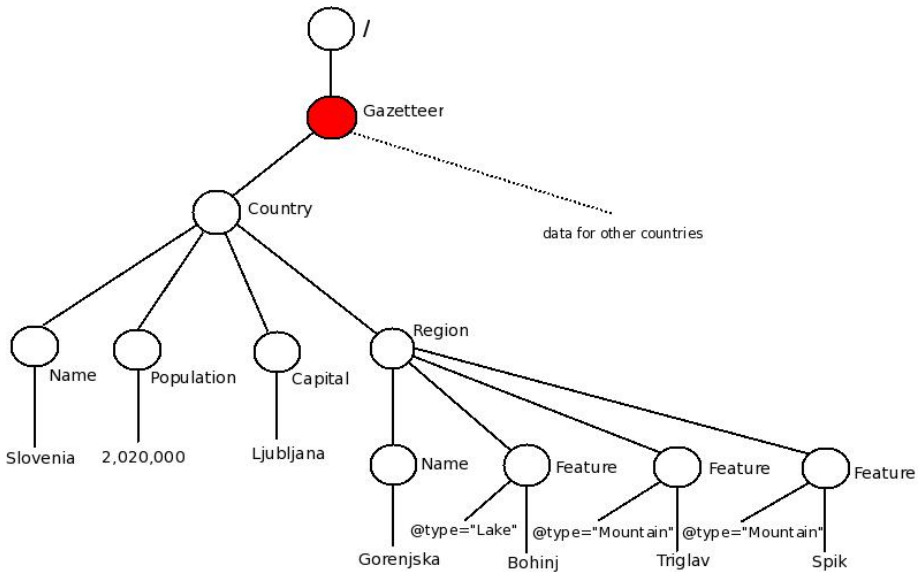
Examples of Location Paths

The next few slides illustrate a selection of location paths applied to the gazetteer example. Each expression appears twice: once using full XPath syntax, and once using a standard abbreviated syntax.

In each case, the nodes identified by the path are highlighted in red, and for a query would be retrieved in document order.

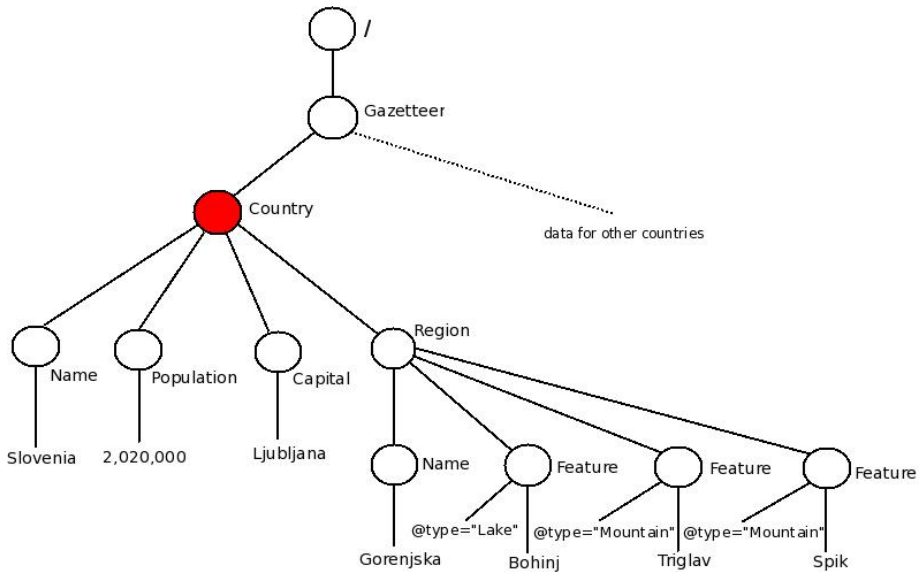
Paths are built up step-by-step as the location path is read from left to right, with a *context node* that travels over the tree according to the components of the location path.

The slash / at the start of a location path indicates that the starting position for the context node is the document root.



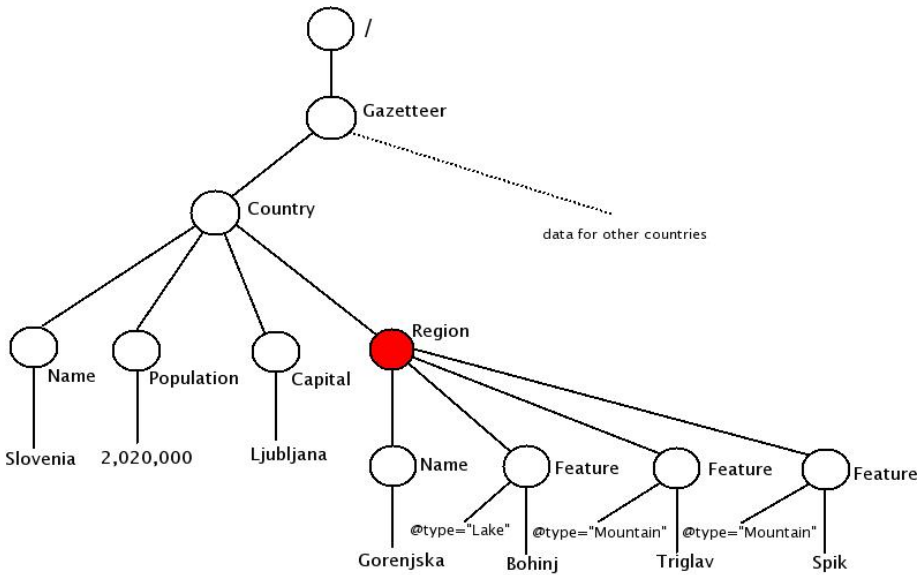
`/child::Gazetteer`

`/Gazetteer`



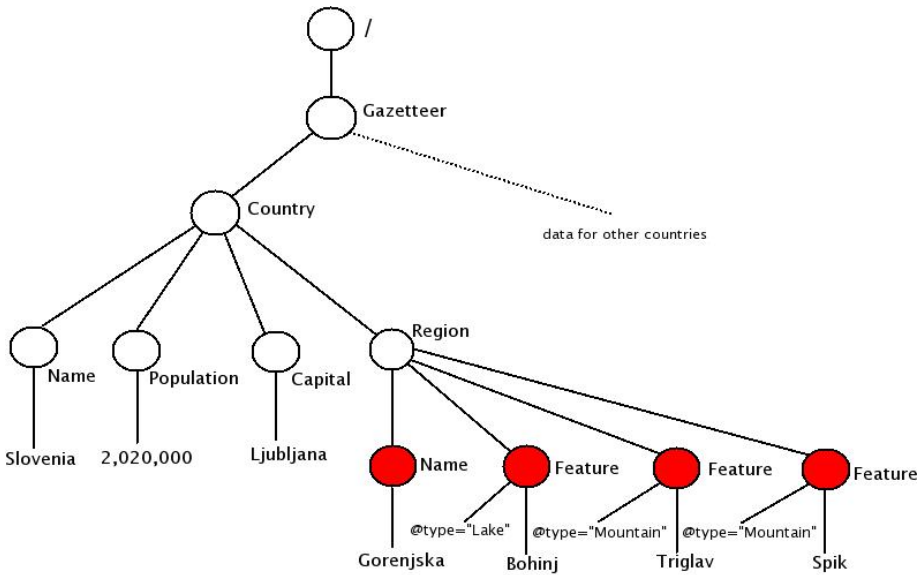
`/child::Gazetteer/child::Country`

`/Gazetteer/Country`



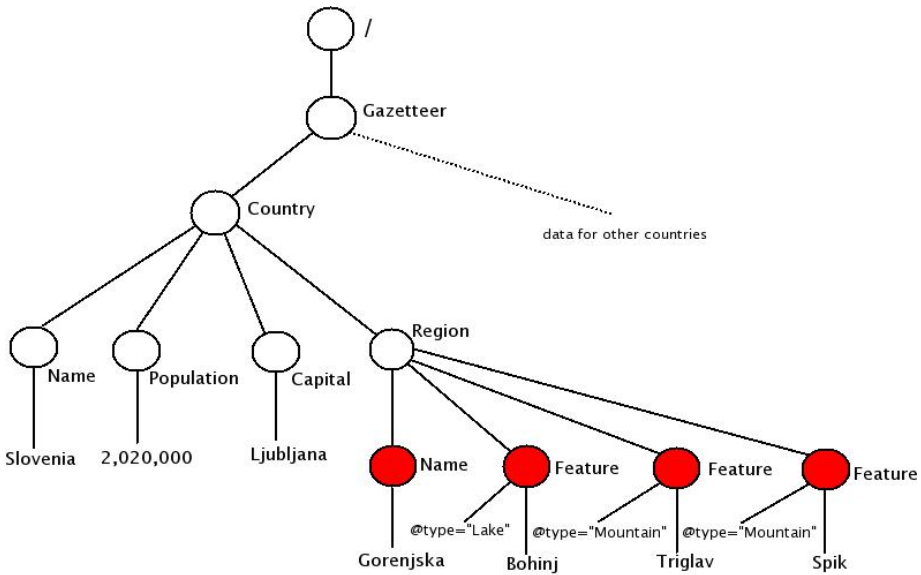
/child::Gazetteer/child::Country/child::Region

/Gazetteer/Country/Region



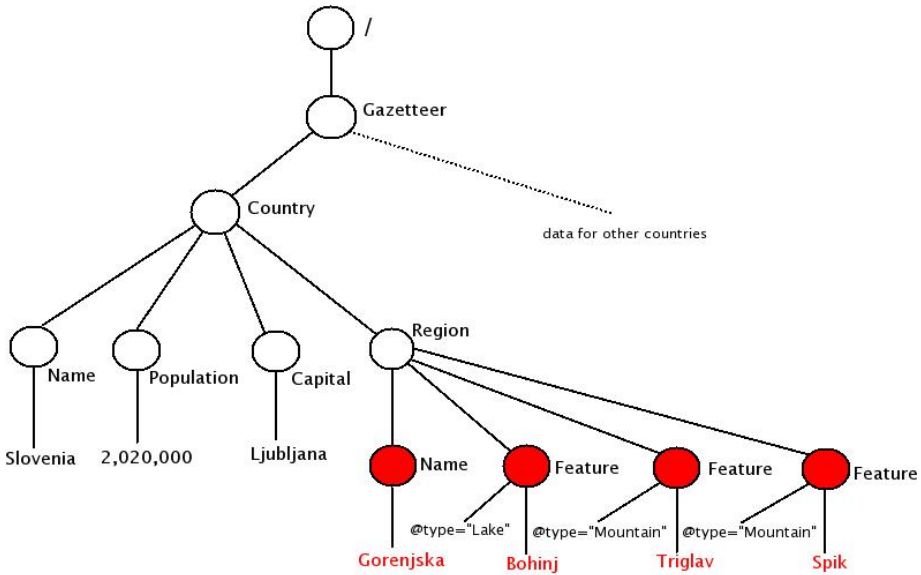
`/descendant::Region/child::*`

`//Region/*`



`/descendant::Region/descendant::*`

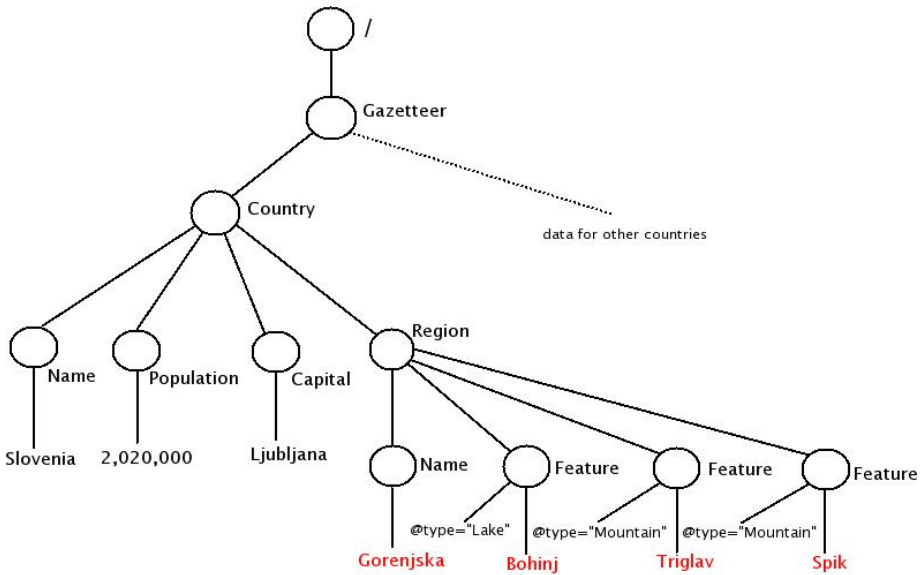
`//Region//*`



```

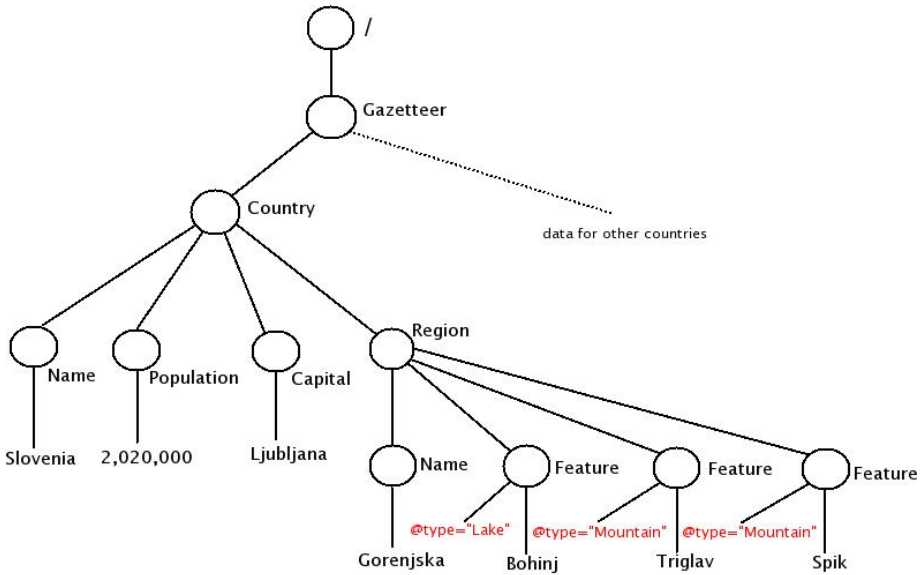
/descendant::Region/descendant::node()
//Region//node()

```

`/descendant::Region/descendant::text()`

`//Region//text()`



`/descendant::Feature/attribute::type`

`//Feature/@type`

Syntax for Location Paths

A *location path* is a sequence of *location steps* separated by a / character.

Each location step has the form

$$\textit{axis} :: \textit{node-test} \textit{predicate}^*$$

- The *axis* indicates which way the context node moves.
- The *node test* selects nodes of an appropriate type.
- The optional *predicates* supply further conditions that need to be satisfied to continue with the path.

The examples so far used the `child` and `descendant` axes; node-tests `node()`, `text()`, `*`, and individual names; and no predicates.

Some Axes

Different axes point in different directions from the current context node.

- **child**: immediate children (attribute nodes don't count)
- **descendant**: any descendants (again, not attribute nodes)
- **parent**: the unique parent (root has no parent)
- **attribute**: all attribute nodes (context node must be an element node)
- **self**: the context node itself
- **descendant-or-self**: the context node together with its descendants.

Some Node Tests

Node tests select among all nodes along the current axis.

- `text()`: nodes with character data.
- `node()`: all nodes.
- `*`: all nodes of the “principal” node type for this axis: for the `attribute` axis, this is attribute nodes; for any other axis, element nodes.
- `name`: element nodes with the given name.

The names used for node tests in the earlier examples were: `Gazetteer`, `Country`, `Region`, `Feature` and `type`.

XPath Abbreviations

Complete path expressions can become cumbersome, and XPath provides a number of abbreviations for the basic operations.

- The `child::` axis is default and can be omitted
- Syntax `@` is an abbreviation for `attribute::`
- Syntax `//` is an abbreviation for `/descendant-or-self::node()/`
- Syntax `..` is an abbreviation for `parent::node()`
- Syntax `.` is an abbreviation for `self::node()`

UTF-8 Encoding

Syntax for Location Paths

A *location path* is a sequence of *location steps* separated by a / character.

Each location step has the form

$$\textit{axis} :: \textit{node-test} \textit{predicate}^*$$

- The *axis* indicates which way the context node moves.
- The *node test* selects nodes of an appropriate type.
- The optional *predicates* supply further conditions that need to be satisfied to continue with the path.

The examples so far used the `child` and `descendant` axes; node-tests `node()`, `text()`, `*`, and individual names; and no predicates.

Some Predicates

The node test in a location step may be followed by zero, one or several *predicates* each given by an expression enclosed in square brackets.

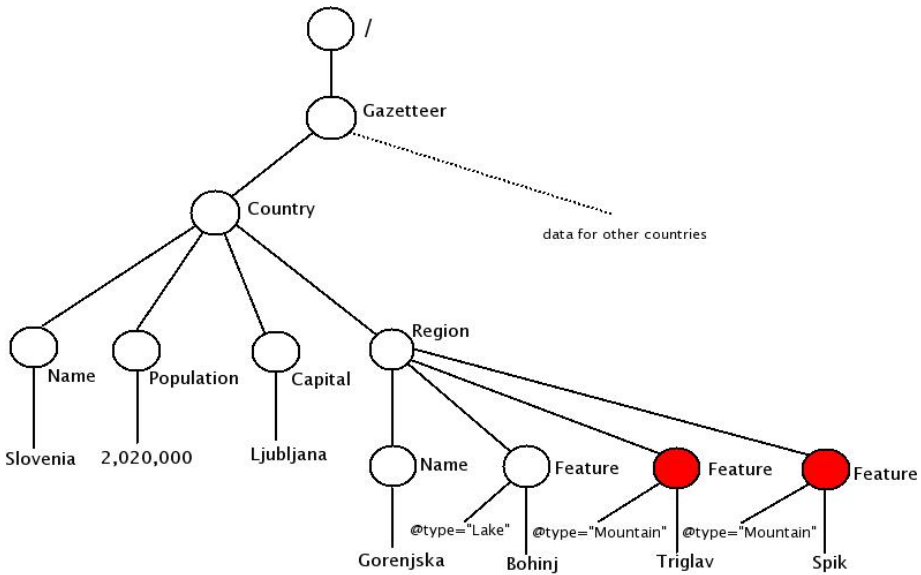
[locationPath]

Selects only those nodes for which there exists a continuation path matching *locationPath*.

[locationPath=value]

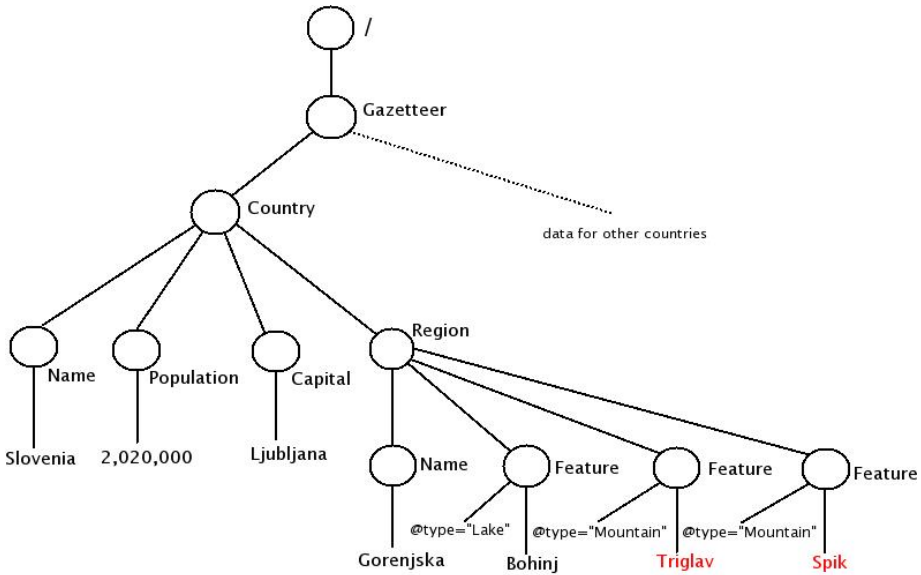
Selects nodes for which there is a continuation path matching *locationPath* where the final node of the path is equal to *value*.

The full syntax of XPath predicate expressions includes arithmetic operations and further path queries, and is beyond the scope of this course.



`/descendant::Feature[attribute::type='Mountain']`

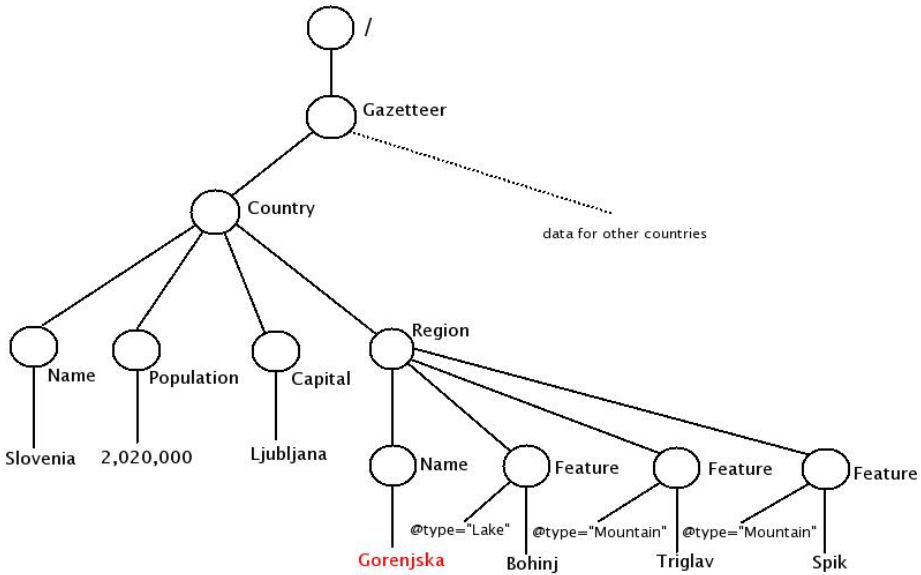
`//Feature[@type='Mountain']`



```

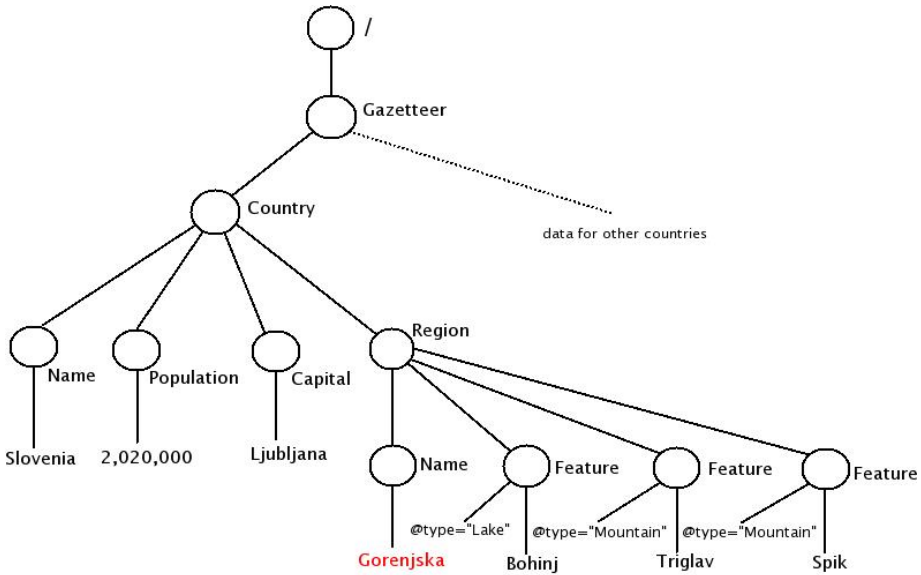
/descendant::Feature[attribute::type='Mountain']/child::text()
//Feature[@type='Mountain']/text()

```



`/descendant::Feature[attribute::type='Mountain']/parent::*/*/child::Name/child::text()`

`//Feature[@type='Mountain']/../Name/text()`



`/descendant::*[Feature/attribute::type='Mountain']/child::Name/child::text()`

`//*[Feature/@type='Mountain']/Name/text()`

XPath as Query Language

These last examples begin to show XPath as a query language, in this case identifying in turn:

- All features which are mountains;
- The names of all mountains;
- The names of all regions containing mountains.

When using XPath in practice, it's often necessary to prefix a location path with a pointer to the relevant XML document:

```
doc("gazetteer.xml")//Feature[@type='Mountain']/text()
```

Subtleties in Complex Queries

Name all countries containing a feature called “Salmon River”

We can select this from a gazetteer with the following XPath expression:

```
//Country[./Feature/text()='Salmon River']/Name/text()
```

Note the use of ‘.’ to start a predicate path at the current context node.

However, this other — apparently very similar — expression won't do:

```
//Country[//Feature/text()='Salmon River']/Name/text()
```

Without ‘.’ the predicate `//Feature/text()` goes back to the root node.

More on XPath

Full XPath has a host of other features, including: navigation based on document order, position and size of context; name spaces; and a rich expression language.

Further Reading

The official W3C specification: <http://www.w3.org/TR/xpath>

Wikipedia on XPath: <https://en.wikipedia.org/wiki/Xpath>

The (wildly optimistic) *10-minute XPath Tutorial*: <http://is.gd/xpath10>

Homework

Tutorial sheet 5 is now online. This involves writing an XML DTD and XPath queries, and running command-line tools which use them.

There's quite a lot to do in this one, so start soon.