

Informatics 1, 2010  
School of Informatics, University of Edinburgh

# **Data and Analysis**

## **Part II**

### **Semistructured Data**

**Alex Simpson**

## Part II — Semistructured Data

XML:

### **II.1 Semistructured data and XML**

### II.2 Structuring XML

### II.3 Navigating XML using XPath

Corpora:

### II.4 Introduction to corpora

### II.5 Querying a corpus

## Recommended reading

[DMS], pp. 227–231, covers the topic, but rather superficially.

For a more in-depth treatment see Chapter 2 of:

[XWT] *An Introduction to XML and Web Technologies*  
A. Møller and M. Schwartzbach  
Addison Wesley, 2006

*“A superb summary of the main Web technologies. It is broad and deep giving you enough detail to get real work done. Eminently readable with excellent examples and touches of humour. This book is a gem.”*

Prof. Philip Wadler, University of Edinburgh

## Background

Relational databases record data in tables conforming to relational schemata. This imposes rigid structure on data

In many situations, it is useful to structure data in a less rigid way; for example:

- when the data needs to be made publicly available in a standard and easily readable data format;
- when we wish to *mark up* (i.e. annotate) existing unstructured data (e.g. text) with additional information (e.g. semantic information);
- when the data possesses a natural hierarchical structure and/or the structure of the data we wish to record varies from item to item.

## Semistructured data

*Semistructured data* imposes a loose structure on data, hence the choice of terminology.

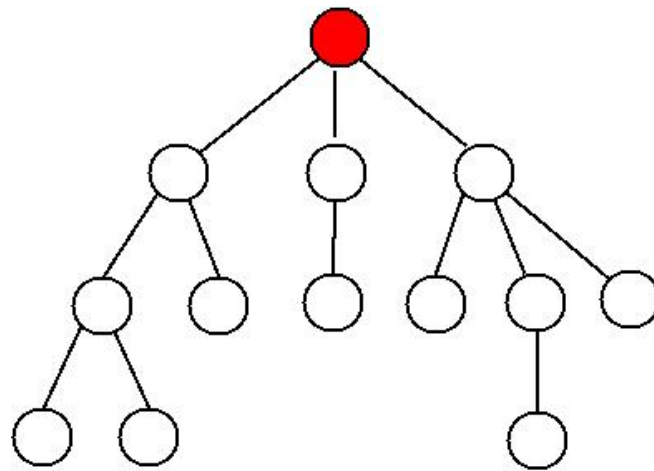
The principal structure imposed on data is that of a *tree*.

Before seeing how trees are used to structure data, we review basic terminology for talking about trees.

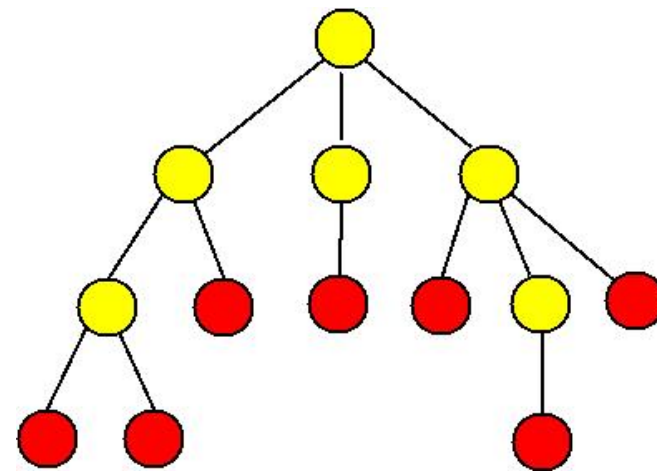
Recall, a *tree* consists of a set of *nodes*, amongst which there is a unique *root node*. For every node in the tree, there is a unique path from the root node to the node.

Nodes separate into two disjoint classes: *leaves* and *internal nodes*.

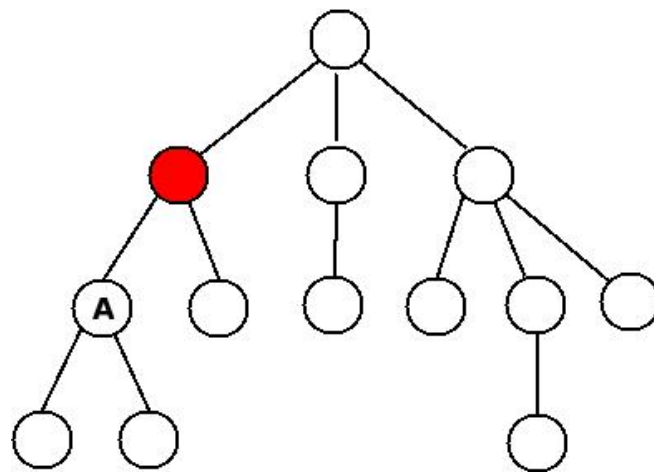
Every node other than the root has a unique *parent* node. Every internal node has a nonempty set of *children* nodes.



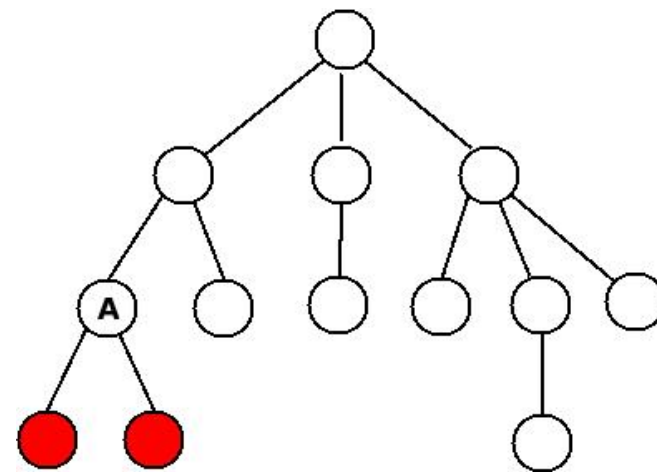
Root node



Leaves and internal nodes



Parent of A



Children of A

## Semistructured data models

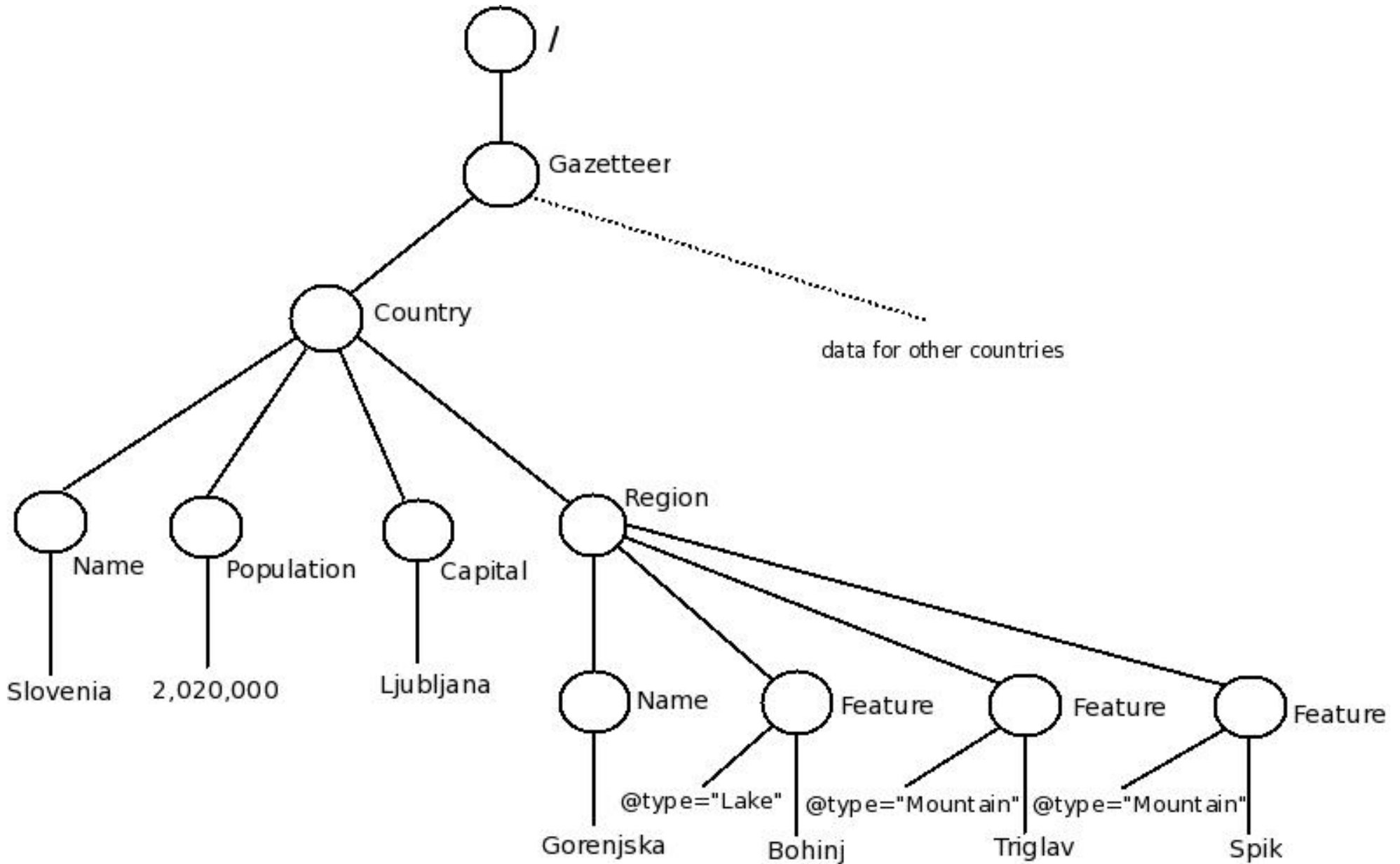
Data is incorporated into a tree structure using a *semistructured data model*.

There are several different such data models.

We shall use the *XPath data model* (chosen because its structure corresponds exactly to XML).

The next slide illustrates an example of data structured according to the XPath data model.

The chosen example, a fragment of a gazetteer, is given because it is one that is naturally accommodated within a hierarchical tree-based structure.





## Types of node in the XPath data model

*Root node.* This is the root of the tree. It is labelled `/`.

*Element nodes.* These are nodes labelled with *element names*, which serve the purpose of categorising the data below them. In the example, the element names are: **Gazetteer**, **Country**, **Name**, **Population**, **Capital**, **Region**, and **Feature**. In the XPath data model, internal nodes other than the root are always element nodes.

The root node is required to have a single element node as child, called the *root element* (since it is root in the tree of all element nodes). In the example, the root element is **Gazetteer**.

*Text nodes.* These are leaves of the tree where textual information is stored. In the example, the text strings "Slovenia", "2,020,000", "Ljubljana", "Gorenjska", "Triglav", "Bohinj" and "Špik" appear at text nodes.

## Attribute nodes

*Attribute nodes* are leaves of the tree in which an *attribute* associated with the parent element node is assigned a value. In the example, we use the @ symbol to identify attributes. There is a single attribute **type**, it is associated with the **Feature** element, and it is assigned the text values "**Lake**" and "**Mountain**".

In the XPath data model, attribute nodes are treated differently from other nodes.

Although the parent of an attribute node is an element node, when we talk about the children of this parent node, attribute nodes are not considered to be amongst them.

Since this can be confusing, explicit warnings will be given in situations in which confusion might arise.

## Understanding the tree

The meaning of the data at a text node depends on the element nodes that appear along the path from the root of the tree to the leaf, and on the values of the attributes to this node.

For example, the path to **Bohinj** is

`/Gazetteer/Country/Region/Feature/`

and the value of the **type** attribute of the associated **Feature** element is "**Lake**". This tells us that Bohinj is a feature in a region in a country in the gazetteer, and that the type of feature is a lake.

Note that to get further information (such as the name of the country, Slovenia), we need to extract it by following another path within the relevant ancestor element (in this case, the **Country** element).

Similarly, the meaning of an element node depends on the path to the node from the root of the tree.

For example, the element **Name** is used in two different ways.

A path **/Gazetteer/Country/Name/** leads to a text node containing the name of a country.

A path **/Gazetteer/Country/Region/Name/** leads to a text node containing the name of a region.

*XML* is a text-based language for presenting exactly the same tree-structured information as the XPath data model.

## Extensible Markup Language (XML)

This is a *markup language*, that is it provides a mechanism, based on *elements* (also called *tags*), for annotating (*marking up*) ordinary text with additional information.

It was developed in the mid 1990's from the Standard General Markup Language (SGML) and Hypertext Markup Language (HTML).

XML has a simple text-based format which provides a convenient basis for making data widely available, e.g. over the web. Indeed, XML has become the *de facto* standard for publishing data on the web.

The next slide presents the gazetteer example in XML format.

The content and structure are identical to that of the tree presented earlier. Only the format is different.

```
<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">Špik</Feature>
    </Region>
  </Country>
  <!-- data for other countries here -->
</Gazetteer>
```

## XML Elements

*Elements* (also called *tags*) are the building blocks of XML documents.

The start of the content of an element *elm* is marked with the *start tag* `<elm>`, and the end of the content is marked with the *end tag* `</elm>`.

Elements must be *properly nested*. Thus,

```
<Country><Region> ... </Region></Country>
```

is legal, whereas

```
<Country><Region> ... </Country></Region>
```

is illegal.

Elements are case sensitive, so **REGION** would be different from **Region**.

The content of the **Capital** element

```
<Capital>Ljubljana</Capital>
```

is the text string "Ljubljana".

The content of the **Region** element consists of one **Name** element together with three **Feature** elements in sequence.

The *root element* **Gazetteer** encloses all information in the document.

Although there are no such examples in the example document, the content of an element may be empty, e.g.,

```
<elm></elm>
```

Such *empty elements* can be abbreviated using a single hybrid tag:

```
<elm/>
```



## Attributes

An element can have descriptive attributes that provide additional information about the element. For example,

```
<Feature type="Mountain"> ... </Feature>
```

sets the attribute **type** of the given **Feature** element to have value **Mountain**.

Note that attribute values are enclosed in quotation marks (either double or single quotes).

It is possible for one element to have several different attributes, with values defined in sequence within the start tag, e.g.

```
<elm attr1="value1" attr2="value2"> ... </elm>
```

## Relating XML and the tree model

The existence of a root element together with the proper nesting of elements ensures that every XML document carries a tree structure in a natural way:

- Each element of the XML document corresponds to an individual element node of the tree.
- The root element of the XML document corresponds to the root element (but *not* the root node) of the tree.
- The text content of an individual XML element corresponds to a child text node of the corresponding element node in the tree.
- An attribute definition in an element's start tag corresponds to a child attribute node of the corresponding element node in the tree.

## Comments and processing instructions

*Comments* can be inserted anywhere in an XML document. Comments start with `<!--` and end with `-->`. They can contain arbitrary text apart from the string `--`.

The full XPath data model also contains *comment nodes* which correspond to XML comments. We do not consider such nodes in our tree model for two reasons:

1. Simplicity.
2. We have included all the types of node that should be used to store data. Comments should instead be used as aids to the interpretation of the data represented.

XML and the XPath data model also allow *processing instructions* to be included. These are beyond the scope of this course.

## Unicode

An XML document is a text document written in *Unicode*.

Unicode is a universal code for “text characters”, currently supporting around 100,000 different characters.

The Unicode characters contain the standard ASCII character set, but also all “characters” in human use worldwide. (The majority of the 100,000 assigned characters are Chinese!)

Each character has an assigned *code point*, which is a number between 0 and 1,114,112.

The actual digital representation of Unicode text depends on a choice of encodings of Unicode character sequences as byte streams. Common choices of encoding are: UTF-8, UTF-16, UTF-32, ISO-8859-1.

## Well-formed documents

An XML document is *well-formed* if it conforms to three guidelines:

- It starts with an XML declaration. (Our example gazetteer document does not!) A suitable such declaration would be:

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

This declares the XML version, and states that UTF-8 character encoding is to be used for Unicode. (Such declarations are not examinable. In Data & Analysis, we are interested in the *content* of a document not in its declaration.)

- It has a root element that contains all other elements.
- All elements are properly nested.

These are minimal requirements on a document. Often there will be other constraints we wish to impose.

## Part II — Semistructured Data

XML:

**II.1** Semistructured data and XML

**II.2 Structuring XML**

**II.3** Navigating XML using XPath

Corpora:

**II.4** Introduction to corpora

**II.5** Querying a corpus

Recommended reading: §§4.1–4.3 of [XWT]

§7.4.2 of [DMS]

## Structuring XML

In a given XML application area, there is often an intended structure that an XML document should possess.

For example, in the **Gazetteer** example, we expect the various elements to respect the natural hierarchy:

- the **Country** elements are inside **Gazetteer**;
- the **Name** (of the country), **Population**, **Capital** and **Region** elements are inside **Country**;
- and the **Name** (of the region) and **Feature** elements are inside **Region**.

Moreover, the **Feature** elements assign a suitable value to the attribute **type**.

## Schema languages for XML

In relational databases, a *schema* specifies the format of a relation (table).

A *schema language* for XML is a language designed for specifying the format of XML documents.

The use of a schema language has two main advantages over giving an informal specification (cf. the informal and partial specification of the **Gazeteer** format on the previous slide):

- It is precise.
- It can be machine checked if an XML document satisfies (*validates*) a given schema specification.

If an XML document  $X$  has the format specified by a given schema  $S$  then we say that  $X$  is *valid* with respect to  $S$ .



## Document Type Definitions

The *Document Type Definition (DTD)* mechanism is a basic schema language for XML.

The language is simple, commonly used, and has been an integrated feature of XML since its inception.

DTD's allow one to specify:

- The elements and entities that can appear in a document.
- What the attributes of the elements are.
- The relationship between different elements including the order of appearance and how they are nested.

We illustrate DTD's by giving an example DTD for a gazetteer format, which validates the XML document on slide II:14.

## Example DTD

```
<!ELEMENT Gazetteer (Country+)>
<!ELEMENT Country (Name,Population,Capital,Region*)>
<!ELEMENT Name (#PCDATA)>
<!ELEMENT Population (#PCDATA)>
<!ELEMENT Capital (#PCDATA)>
<!ELEMENT Region (Name,Feature*)>
<!ELEMENT Feature (#PCDATA)>
<!ATTLIST Feature type CDATA #REQUIRED>
```

## Understanding the example DTD

```
<!ELEMENT Gazetteer (Country+)>
```

This states that the **Gazetteer** element consists of one or more **Country** elements.

```
<!ELEMENT Country (Name,Population,Capital,Region*)>
```

This states that a **Country** element consists of: one **Name** element, followed by one **Population** element, followed by one **Capital** element, followed by zero or more **Region** elements.

```
<!ELEMENT Name (#PCDATA)>
```

This states that the **Name** element contains text. The keyword **#PCDATA** abbreviates “parsed character data”.

```
<!ELEMENT Region (Name,Feature*)>
```

This states that a **Region** element consists of: one **Name**, followed by zero or more **Feature** elements.

```
<!ELEMENT Feature (#PCDATA)>
```

This states that the **Feature** element has text content.

```
<!ATTLIST Feature type CDATA #REQUIRED>
```

This states that the **Feature** element has an attribute **type**, and that the value of the attribute should be a text string (**CDATA** abbreviates “character data”). Moreover, it is *required* that every **Feature** element in the document must assign a value to the **type** attribute.

## General format of element declarations

An *element declaration* has the structure:

```
<!ELEMENT elementName (contentType)>
```

There are four possible content types:

1. **EMPTY** indicating that the element has no content, i.e. it is an *empty element* as defined on slide II:16.

2. **ANY** indicating that any content is permitted.

Nevertheless elements that appear within the element content must themselves be declared by corresponding element declarations.

3. **#PCDATA** indicating text content.

(In fact this is an instance of a more general *mixed content* format, which we shall not consider further.)

4. A *regular expression* of element names.

Regular expressions were introduced in Inf1 Computation and Logic.

DTD's make use of the following format for regular expressions.

- Any element name is a regular expression.  
(The element names are the *alphabet* for the regular expressions.)
- ***exp1, exp2***: first ***exp1*** then ***exp2*** in sequence.
- ***exp\****: zero or more occurrences of ***exp***.
- ***exp?***: zero or one occurrences of ***exp***.
- ***exp+***: one or more occurrences of ***exp***.
- ***exp1 | exp2***: either ***exp1*** or ***exp2***.

## General format of attribute declarations

The attributes of an element are declared separately to the element declaration. The general format is:

```
<!ATTLIST elementName (attName attType default)+>
```

This declares a list of at least one attribute for the element *elementName*.

For each entry in the list:

- *attName* is the attribute name
- *attType* is a type for the value of the attribute.
- *default* specifies whether the attribute is required or optional, and may specify a default value for the attribute.

We shall consider only the following attribute types:

- *String type*: **CDATA** means that the attribute may have any text string as its value.
- *Enumerated type*:  $(s_1 \mid s_2 \mid \dots \mid s_n)$  means that the attribute must take one of the strings  $s_1, s_2, \dots, s_n$  as its value.

And the following default options.

- *Required*: **#REQUIRED** means that the attribute must be explicitly assigned a value in every start tag for the element.
- *Optional*: **#IMPLIED** means it is optional whether a value is assigned to the attribute or not.
- *Default*: A fixed string can be specified as the default value for the attribute to take if no explicit value is given in the element's start tag.



## A variation on the example

Consider replacing the attribute declaration in the example DTD with the following declaration.

```
<!ATTLIST Feature type (Mountain|Lake|River) "Mountain">
```

With this new (but not with the original) declaration:

```
<Feature>Ben Nevis</Feature>
```

would be a valid **Feature** element. The **type** attribute would be given the default (and correct) default value **Mountain**.

The element below is not valid with respect to the new DTD (although it is valid for the original DTD)

```
<Feature type="Castle">Eilean Donan</Feature>
```

because **Castle** is not one of the specified values for **type**.

## Document type declaration

A *document type declaration* can appear in an XML document between the XML declaration and the root element. It links the XML document to a DTD schema intended to specify the structure of the document.

The usual format of a document type declaration is:

```
<!DOCTYPE rootName SYSTEM "URI">
```

where *rootName* is the name of the root element, and *URI* is the *Uniform Resource Indicator* of the intended DTD.

An alternative (illustrated on the next slide) is to include the DTD within the XML document itself, using an *internal declaration*

```
<!DOCTYPE rootName [DTD]>
```

## Example internal document type declaration

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE Gazetteer [
<!ELEMENT Gazetteer (Country+)>
<!ELEMENT Country (Name,Population,Capital,Region*)>
<!ELEMENT Name (#PCDATA)>
<!ELEMENT Population (#PCDATA)>
<!ELEMENT Capital (#PCDATA)>
<!ELEMENT Region (Name,Feature*)>
<!ELEMENT Feature (#PCDATA)>
<!ATTLIST Feature type CDATA #REQUIRED>
]>
<Gazetteer>...</Gazetteer>
```

## Limitations of DTD's

One of the strengths of the DTD mechanism is its essential simplicity.

However, it is inexpressive in several important ways, and this severely limits its usefulness. For example, three weaknesses are:

- Elements and attributes cannot be assigned useful types.
- It is impossible to place constraints on data values.
- There are restrictions on how character data and elements can be combined (they can only be combined as *mixed content*), and there are also undesirable technical restrictions on the forms of regular expression allowed when declaring the structure of elements.

These issues and others have been dealt with through the development of more powerful, but more complex, XML format languages, such as XML Schema (which lie beyond the scope of Data & Analysis.)

## Publishing relational data as XML

A common application of XML is as a format for publishing data from relational databases.

The benefit of XML for this is that its simple text format makes the data easily readable and transferable across platforms.

The generality and flexibility of the XML format means that there are many ways to translate relational data into XML.

We illustrate one simple approach using example data from previous lectures (cf. slide I:99).

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE UniversityData [
  <!ELEMENT UniversityData (Students,Courses,Takes)>
  <!ELEMENT Students (Student*)>
  <!ELEMENT Student (mn,name,age,email)>
  <!ELEMENT Courses (C*)>
  <!ELEMENT C (code,name,year)>
  <!ELEMENT Takes (T*)>
  <!ELEMENT T (mn,name,mark)>
  <!ELEMENT mn (#PCDATA)>
  <!ELEMENT name (#PCDATA)>
  <!ELEMENT age (#PCDATA)>
  <!ELEMENT email (#PCDATA)>
  <!ELEMENT code (#PCDATA)>
  <!ELEMENT year (#PCDATA)>
  <!ELEMENT mark (#PCDATA)>
]>
```

```
<UniversityData>
<Students>
  <Student> <mn>s0456782</mn> <name>John</name>
    <age>18</age> <email>john@inf</email> </Student>
  <Student> <mn>s0412375</mn> <name>Mary</name>
    <age>18</age> <email>mary@inf</email> </Student>
  <Student> <mn>s0378435</mn> <name>Helen</name>
    <age>20</age> <email>helen@phys</email> </Student>
  <Student> <mn>s0189034</mn> <name>Peter</name>
    <age>22</age> <email>peter@math</email> </Student>
</Students>
<Courses>
  <C><code>inf1</code><name>Informatics 1</name><year>1</year></C>
  <C><code>math1</code><name>Mathematics 1</name><year>1</year></C>
</Courses>
<Takes>
  <T><mn>s0412375</mn><code>inf1</code><mark>80</mark></T>
  <T><mn>s0378435</mn><code>math1</code><mark>70</mark></T>
</Takes>
</UniversityData>
```

## Efficiency

Relational database systems are optimised for storage efficiency.

As we have seen, the XML version of relational data is extremely verbose.

Nevertheless, XML can still be stored efficiently using *data compression* (which can be optimised for XML).

Furthermore, once published XML data has been downloaded, it can be converted back to relational data so it can be stored efficiently in a local database system.

Converting XML to back to relational data has the benefit of enabling the data to be queried using relational database technology (i.e., SQL).

An interesting alternative is to apply newer technology for directly querying XML.



## Part II — Semistructured Data

### XML:

#### II.1 Semistructured data and XML

#### II.2 Structuring XML

#### **II.3 Navigating XML using XPath**

### Corpora:

#### II.4 Introduction to corpora

#### II.5 Querying a corpus

### Recommended reading:

§§3.1–3.4 of [XWT]

pp. 948–949 of [DMS] (superficial coverage only)

On-line XPath tutorial: <http://www.w3schools.com/xpath/>

## How do we extract data from an XML document?

Since an XML document is a text document, one option is to use methods based on text search.

But this ignores the element structure of the document.

A better alternative is to use a dedicated language for forming queries based on the *tree structure* of an XML document

This has many uses, for example:

- Performing relational-database-type queries directly on data published as XML
- Extracting annotated content from marked-up text documents
- All queries that exploit the tree structure of XML

## XQuery and XPath

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

However, the XQuery language is too complex for this course. (See [XWT] for further information.)

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

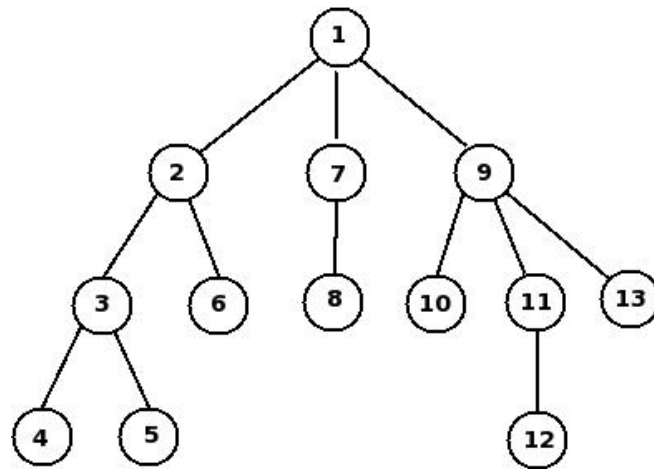
XPath can be viewed as a rudimentary query language in its own right.

It is also an important component of many XML application languages other than XQuery (e.g., XML Schema, XSLT, XLink, XPointer).

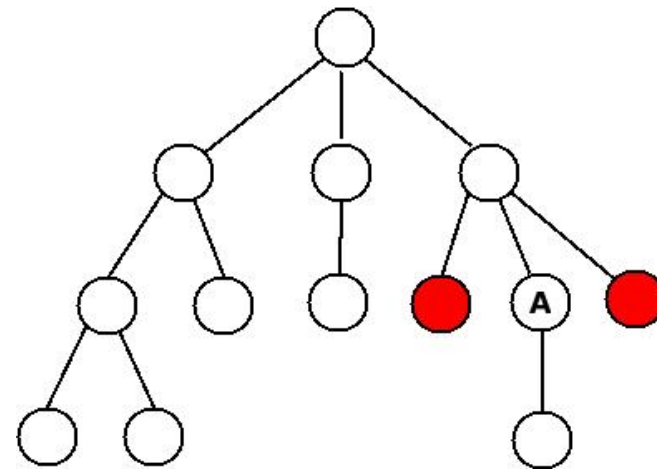
## Location paths

A *location path* (a.k.a. *path expression*) retrieves a *set* of nodes from an XML document tree.

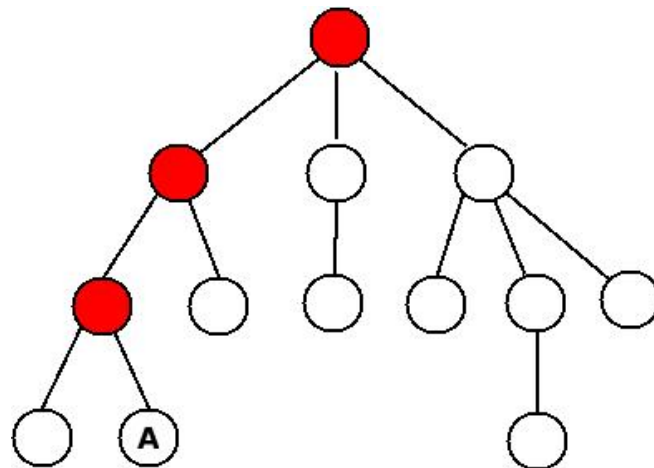
- The location path describes a set of possible paths from the root of the tree.
- The set of nodes retrieved is the set of all nodes reached as final destinations of the described paths.
- This set of nodes is returned as a list of nodes (without duplicates) sorted in *document order* (the order in which the nodes appear in the XML document)



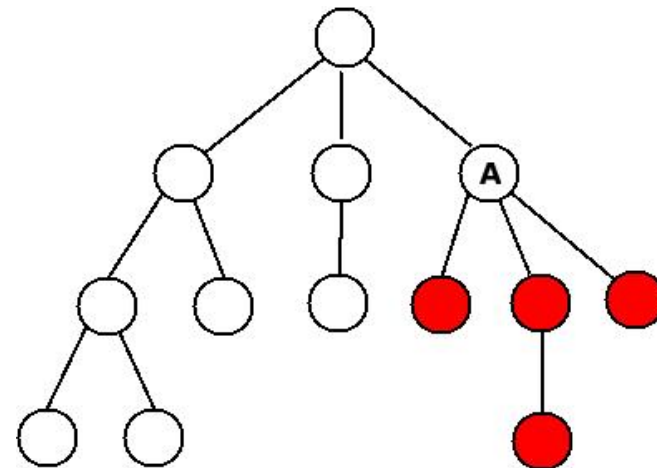
Document order



Siblings of A



Ancestors of A



Descendants of A

## Example location paths

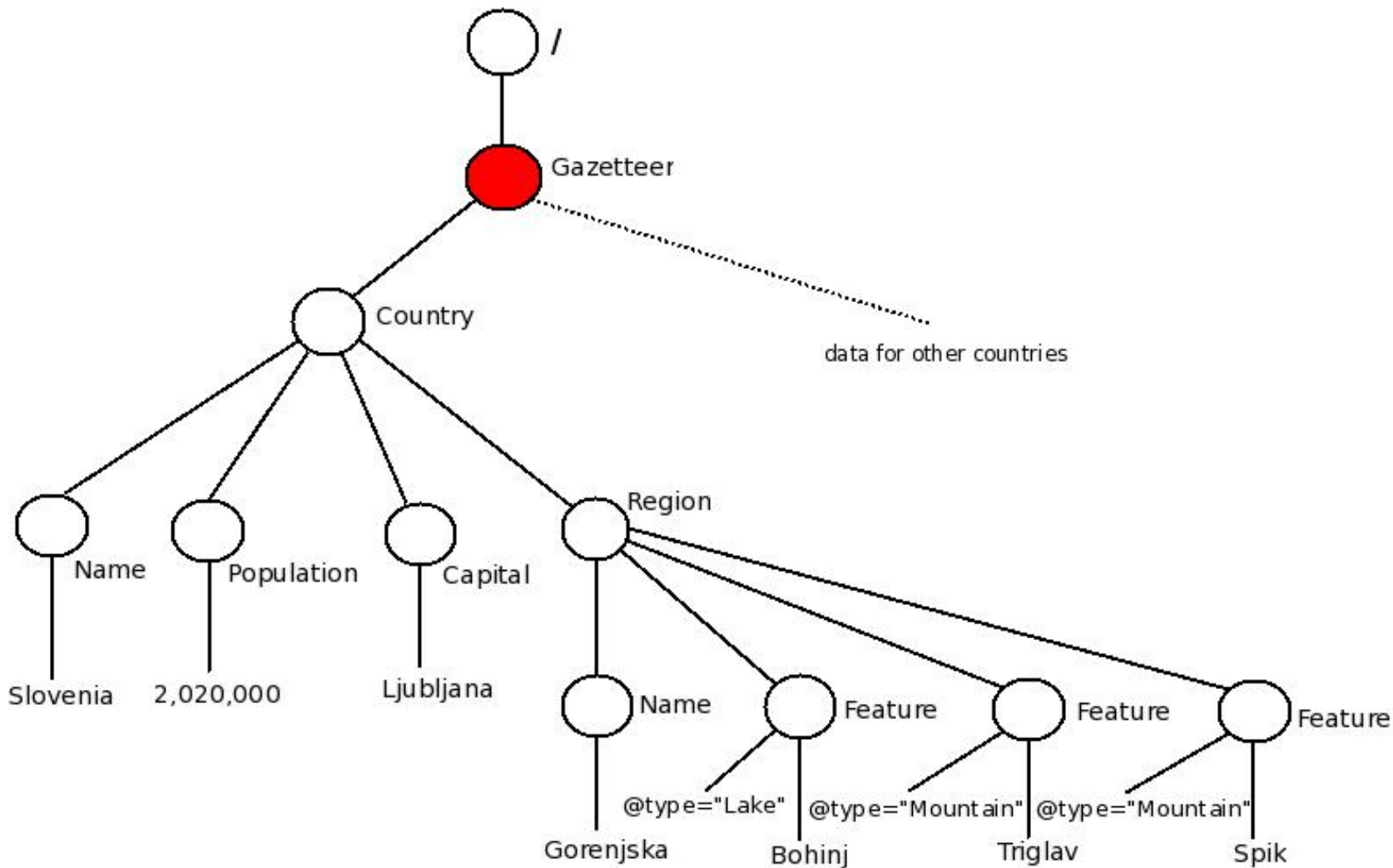
The next few slides illustrate a selection of location paths. Each is given twice: above using the full XPath syntax, and below using a convenient abbreviated syntax.

In each case, the retrieved nodes are highlighted in red. These nodes will be returned as a list in document order.

Paths are built up step-by-step as the location path is read from left-to-right.

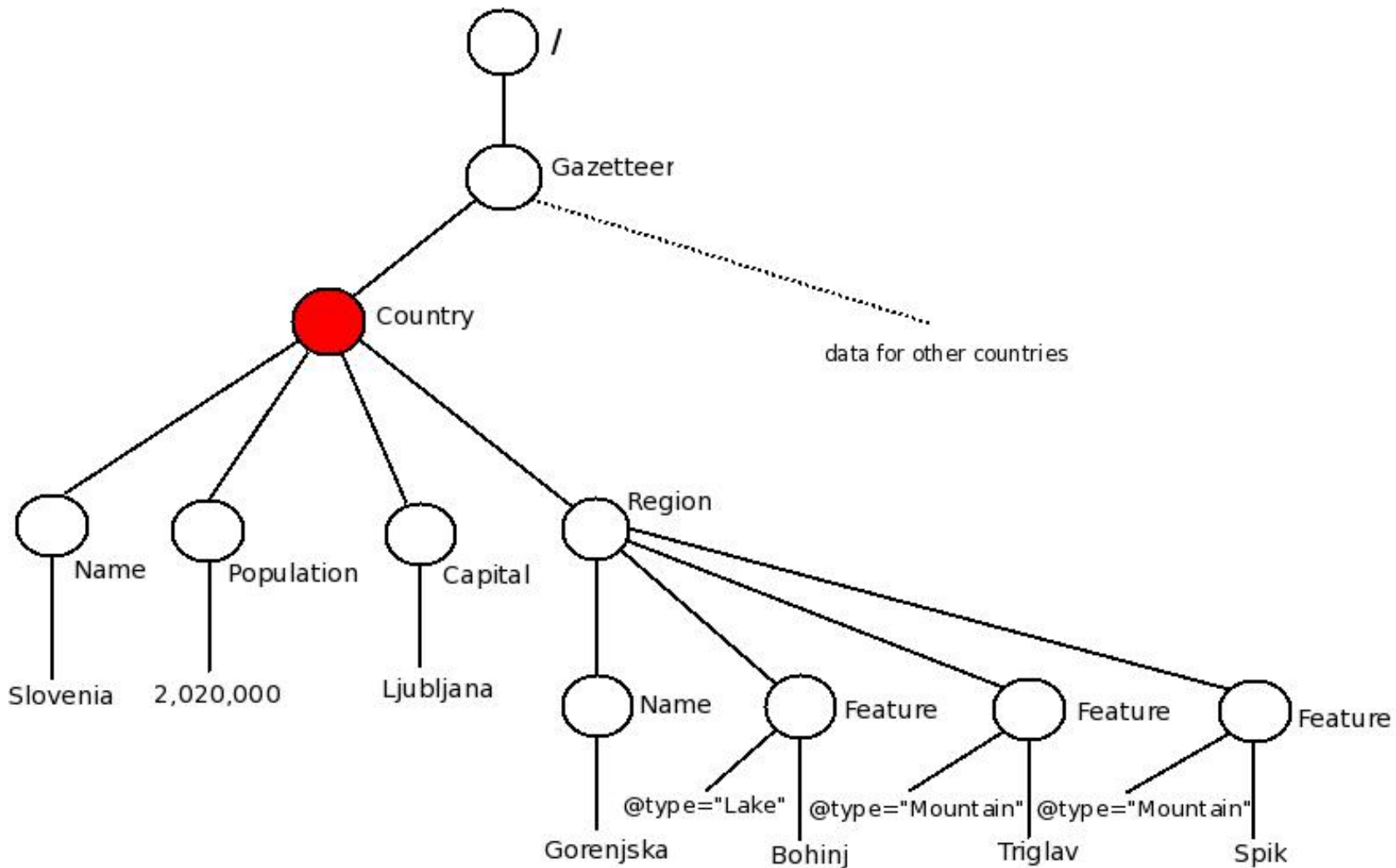
Each path is constructed by a *context node* that travels over the tree, according to certain rules, depending on the continuation of the location path expression.

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



`/child::Gazetteer`

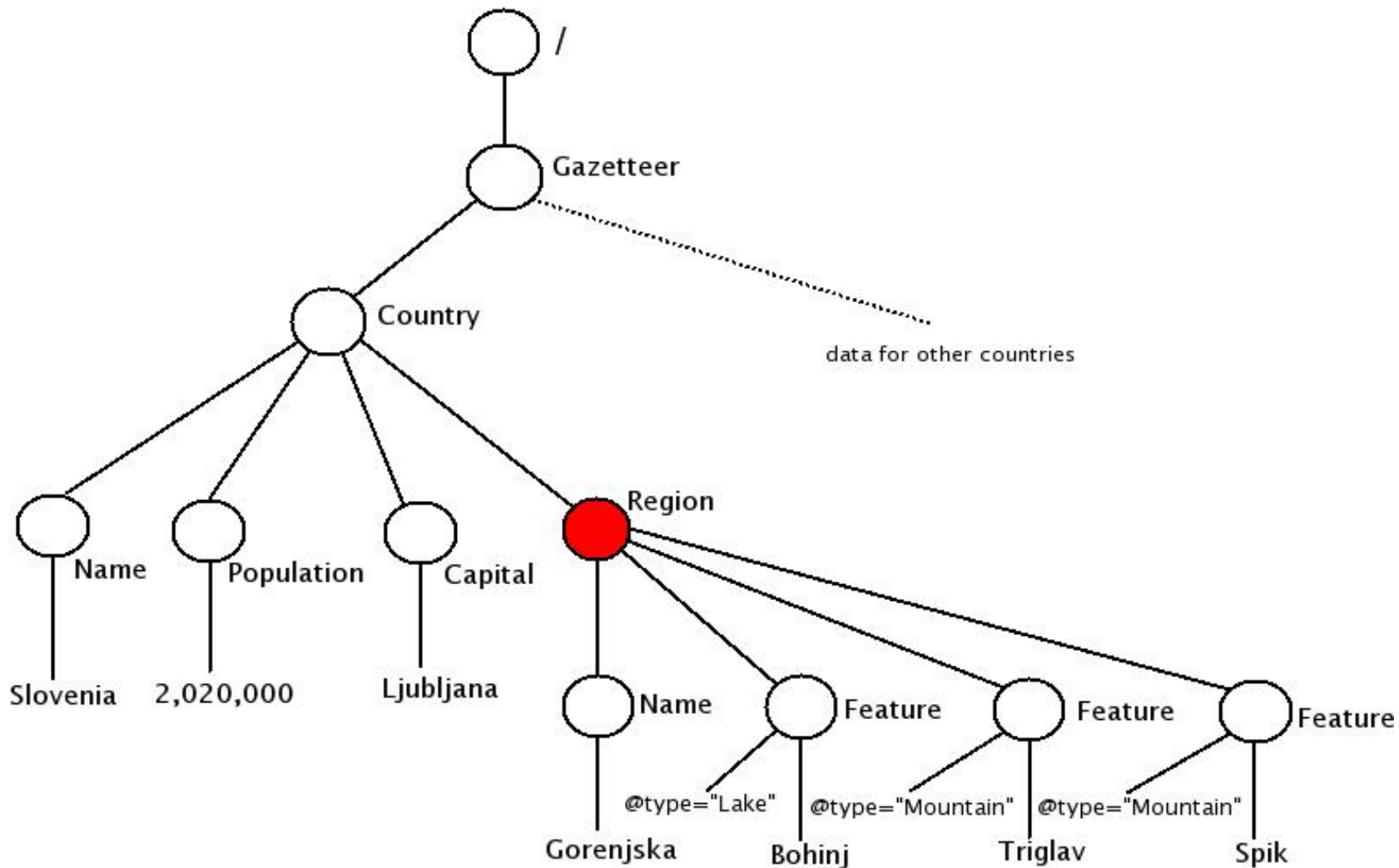
`/Gazetteer`



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

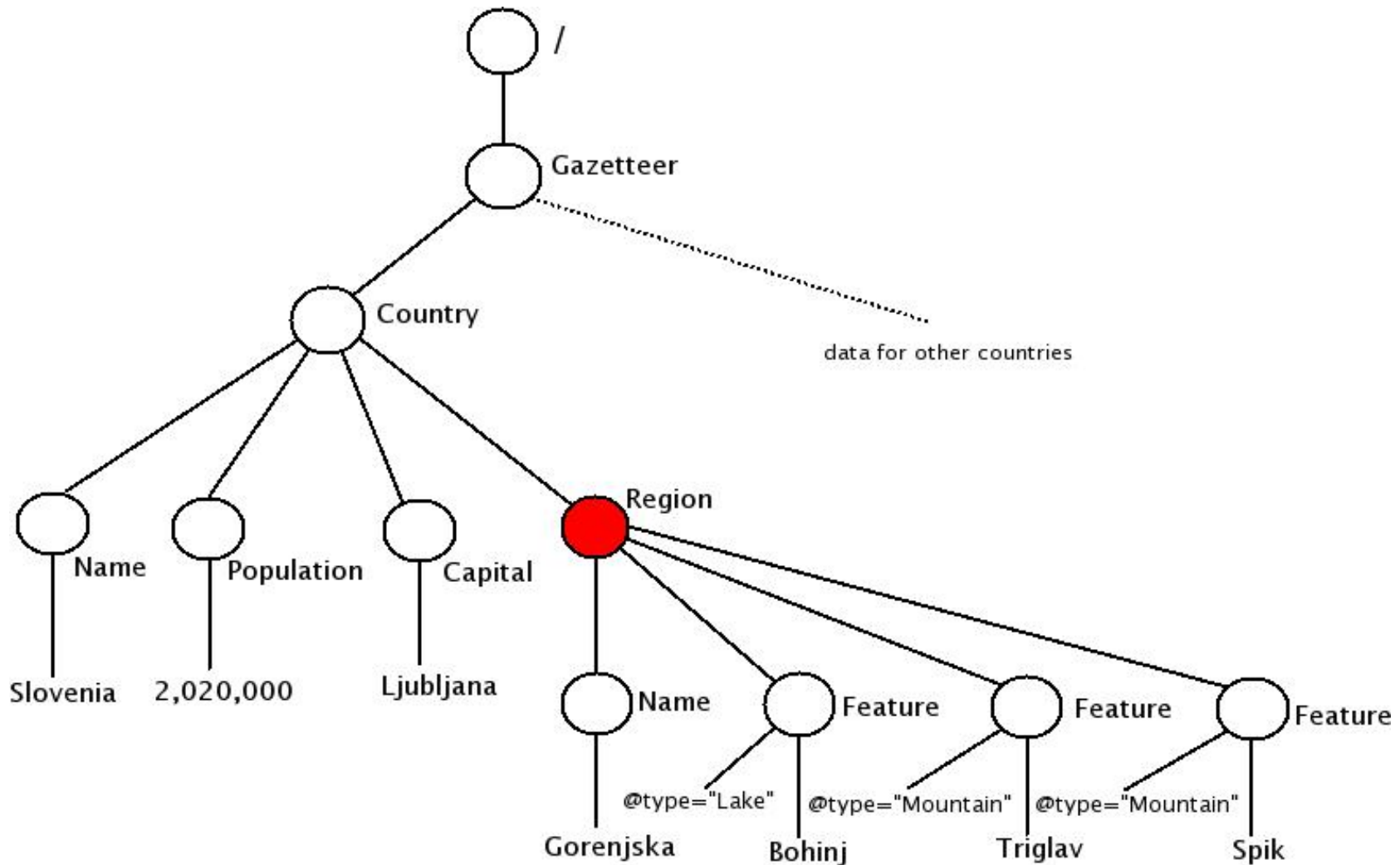
**`/Gazetteer/Country`**





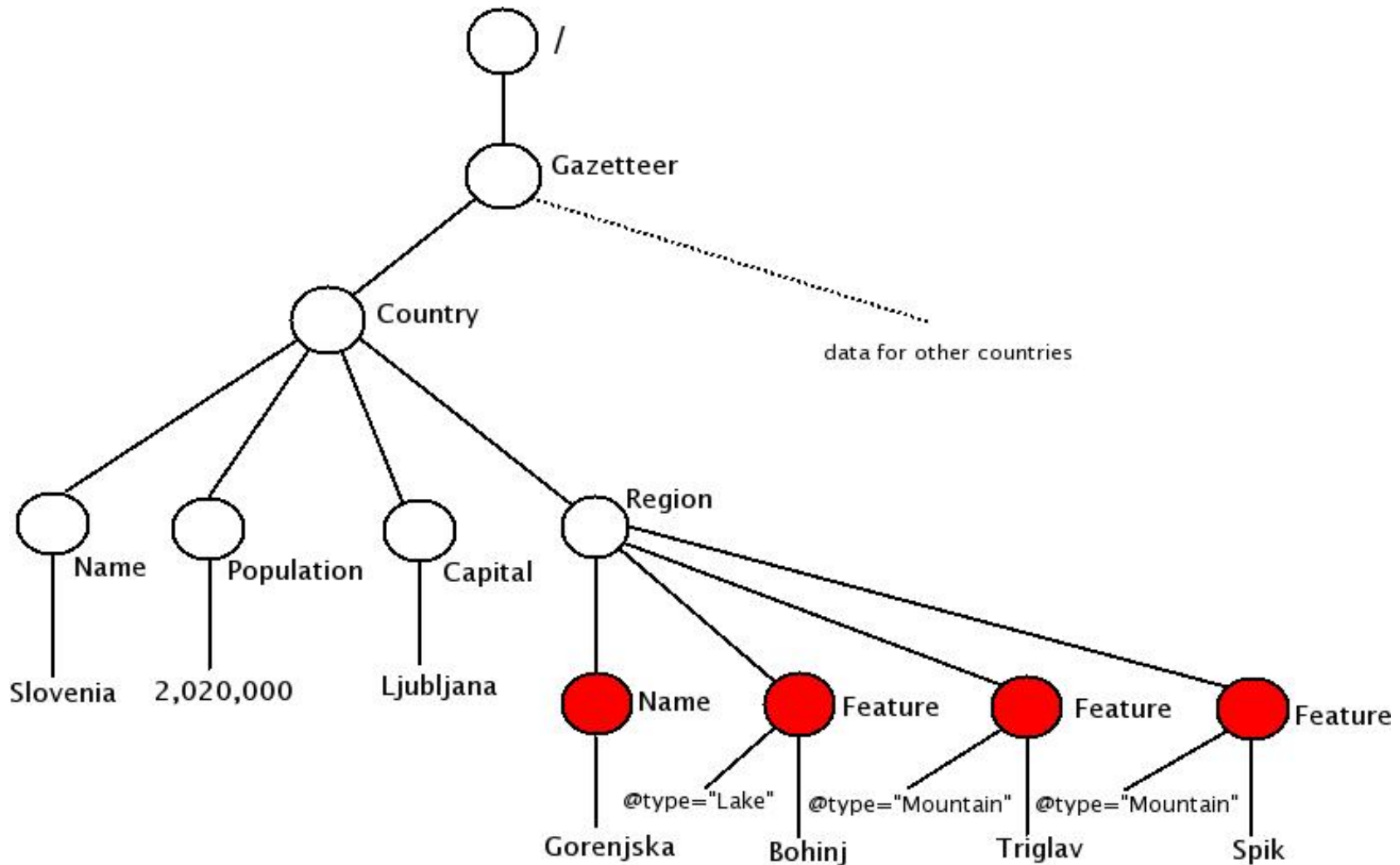
`/child::Gazetteer/child::Country/child::Region`

`/Gazetteer/Country/Region`



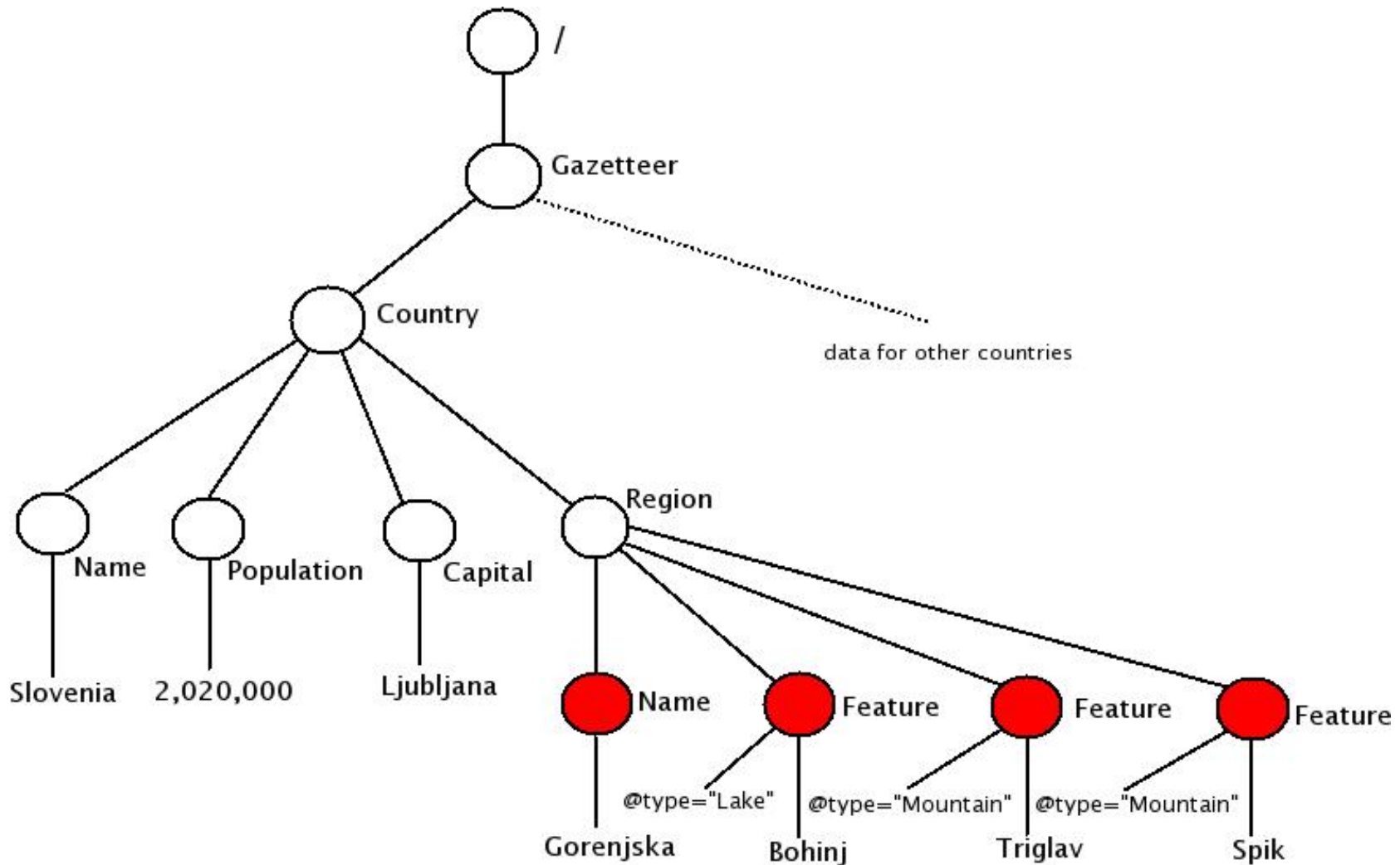
`/descendant::Region`

`//Region`



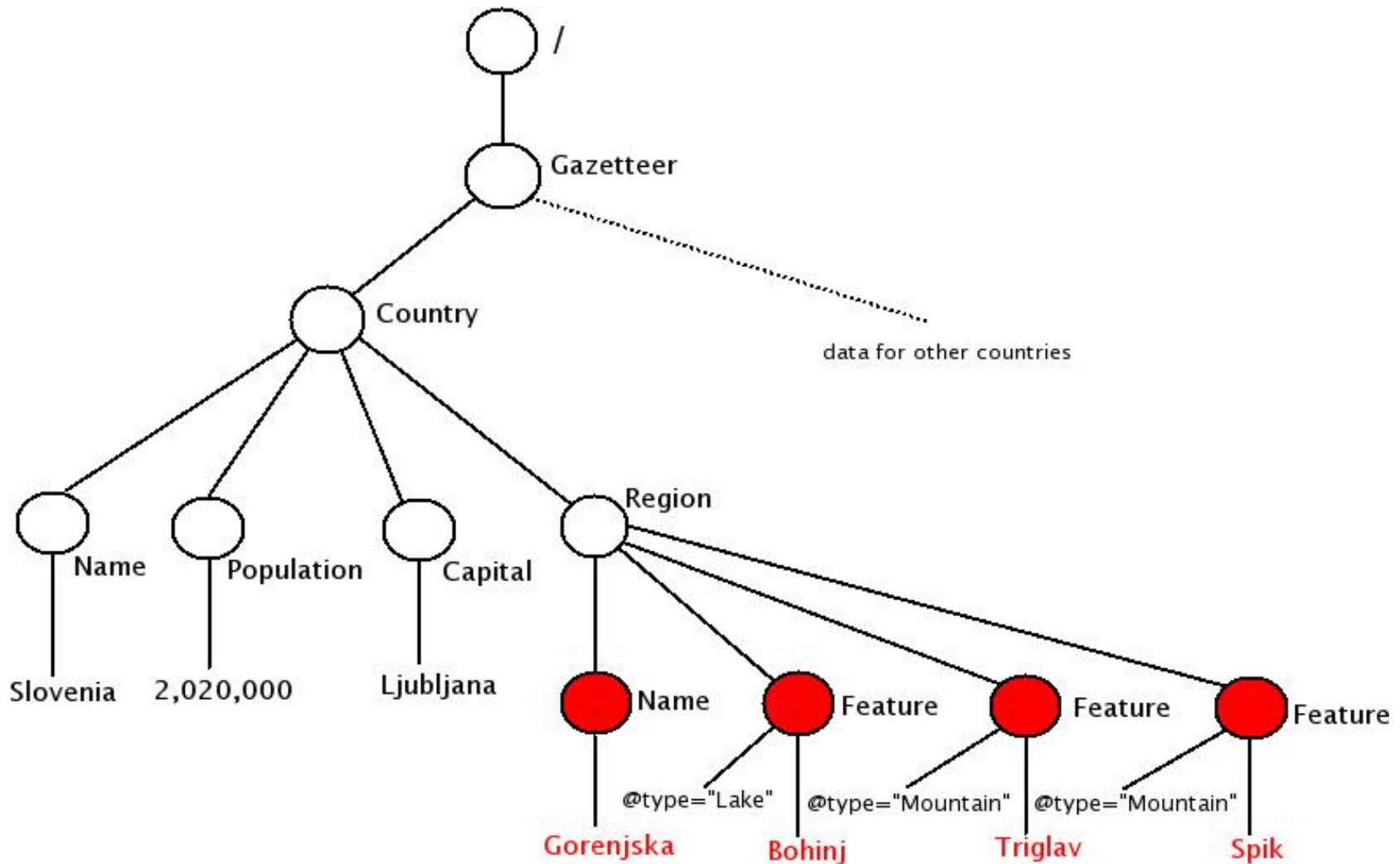
`/descendant::Region/child::*`

`//Region/*`



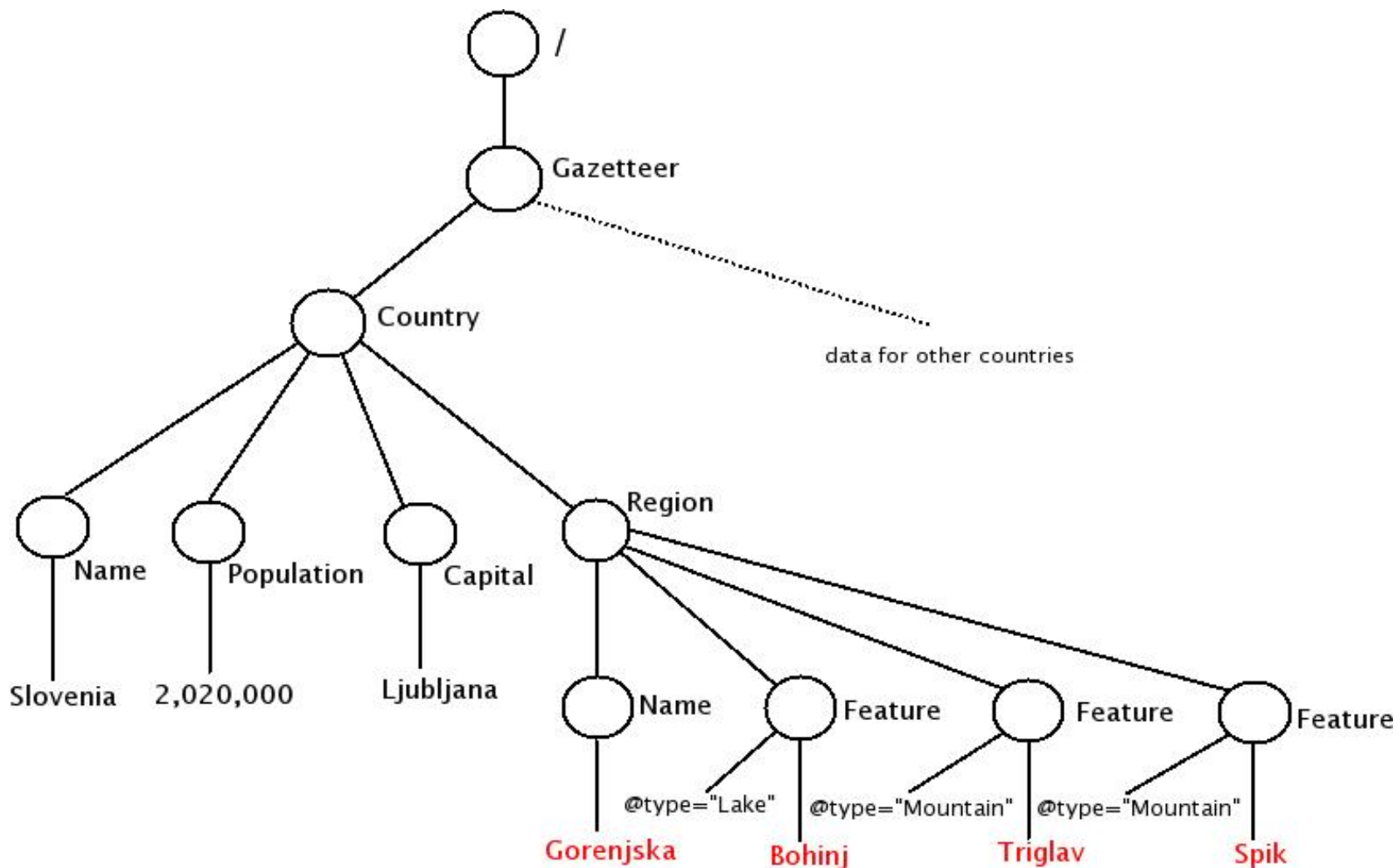
`/descendant::Region/descendant::*`

`//Region//*`



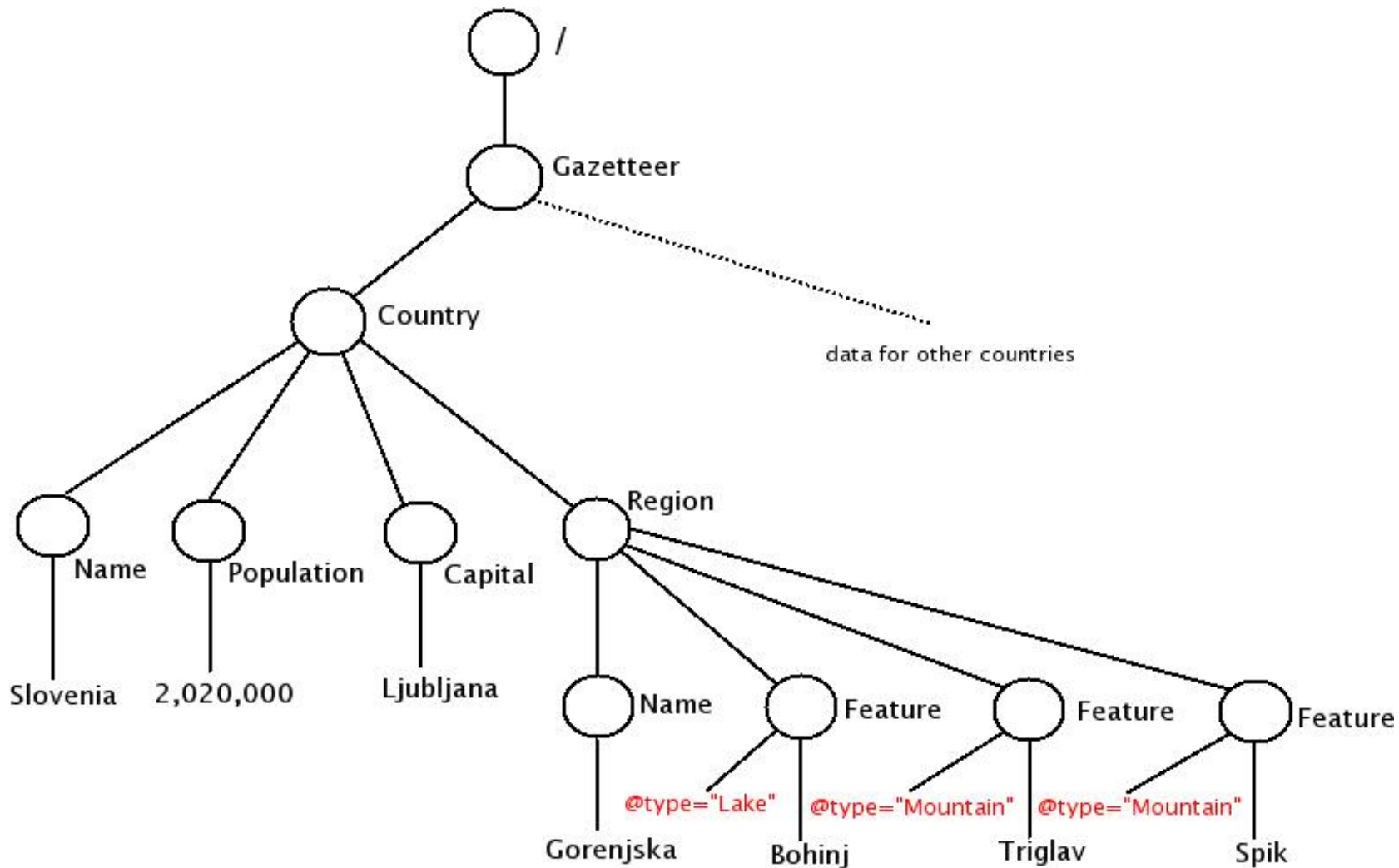
```
/descendant::Region/descendant::node()
```

```
//Region//node()
```



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

```
//Region//text()
```



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

`//Feature/@type`

## General unabbreviated syntax of location paths

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

A *location step* has the form

***axis::nodeTest predicate\****

- The *axis* tells the context node which way to move.
- The *node test* selects nodes of an appropriate type from the tree.
- The optional *predicates* supply conditions that need to be satisfied for the path to be allowed to count towards the result.

N.B., the previous examples contained only axes and node tests.



## A selection of axes

- **child**: the children of the context node (remember, an attribute node does not count as a child node)
- **descendant**: the descendants of the context node (again, an attribute node does not count as a descendant).
- **parent**: the unique parent of the context node (where the context node must not be the root node).
- **attribute**: all attribute nodes of the context node (which must be an element node).
- **self**: the context node itself (this is useful in connection with abbreviations).
- **descendant-or-self**: the context node together with its descendants.

## A selection of node tests

Node tests filter the nodes selected by the current axis according to the type of node.

- **text ()** : selects only character data nodes.
- **node ()** : selects all nodes.
- **\*** : if the axis is **attribute** then all attribute nodes are selected; for any other axis, all element nodes are selected.
- **name** : selects the nodes with the given name.

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

## 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.

Common examples of predicates are:

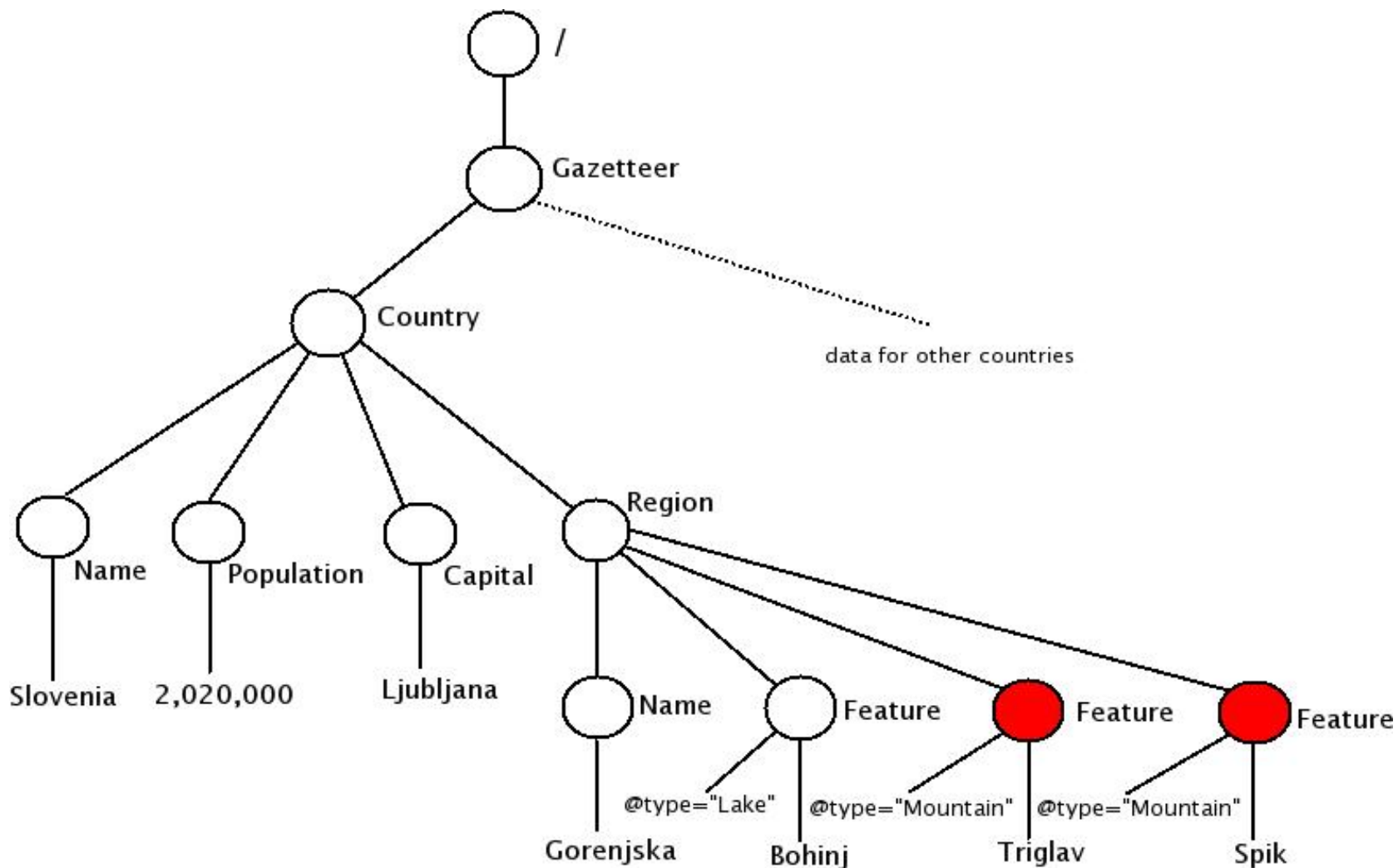
- **[ *locationPath* ]**

This selects only those nodes for which there exists a continuation path (from the current node) matching *locationPath*.

- **[ *locationPath = value* ]**

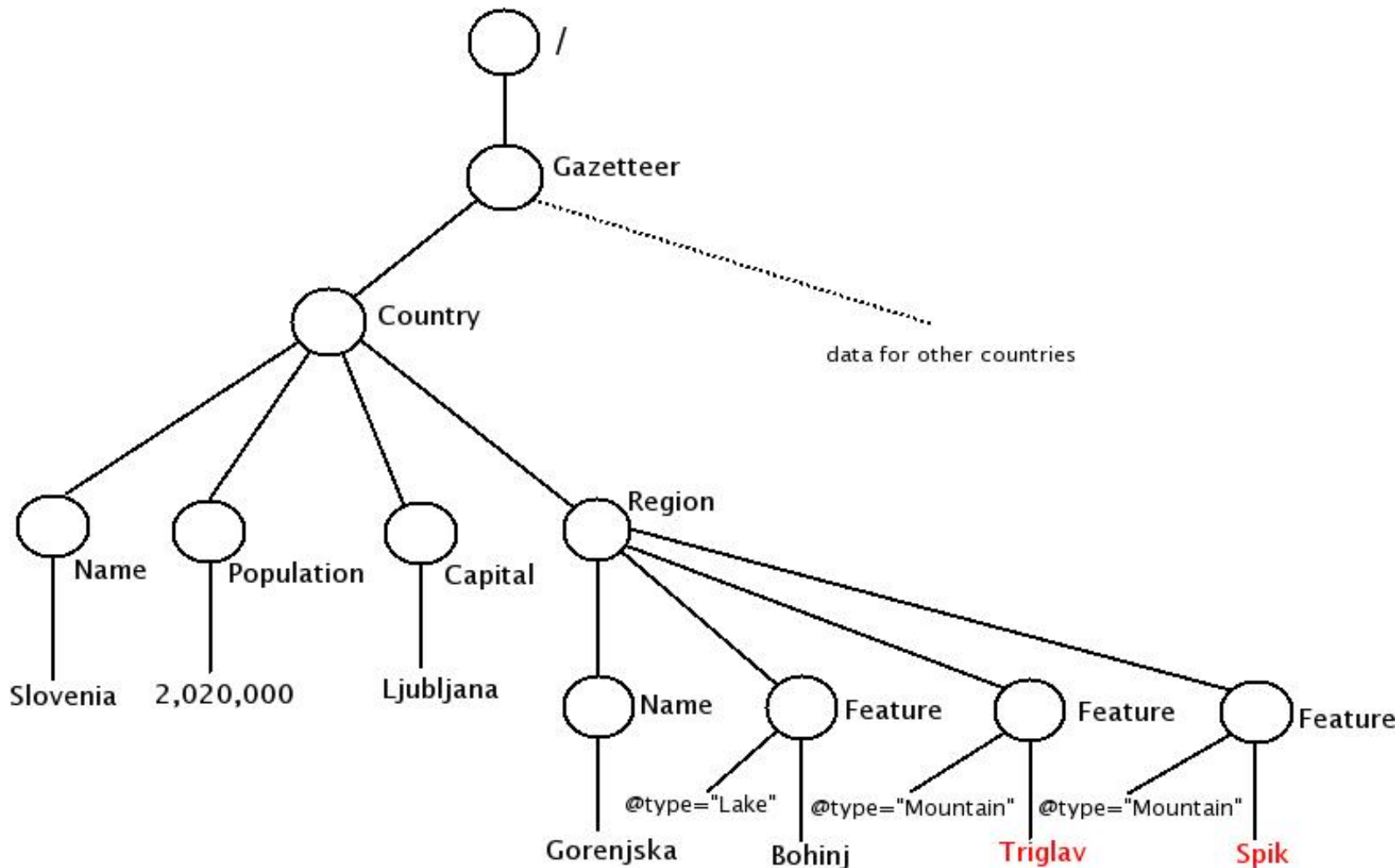
Selects those nodes for which there exists a continuation path matching *locationPath* such that the final node of the path is equal to *value*.

The full syntax of XPath predicate expressions is rather powerful, but beyond the scope of the course.



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

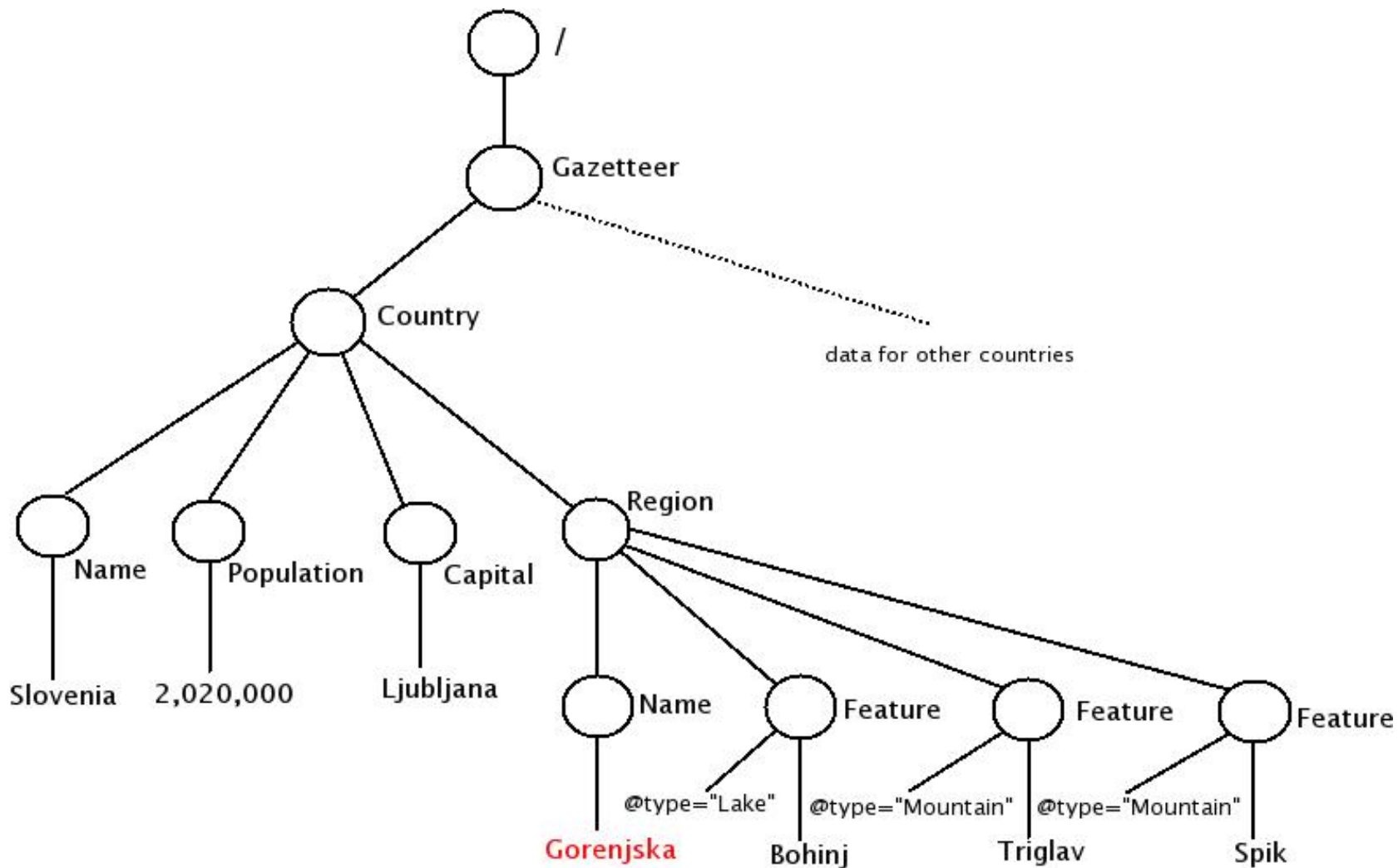
```
//Feature[@type='Mountain']
```



```

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

```



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

## XPath as a query language

The previous examples illustrate XPath as a rudimentary query language.

The queries formulated are:

- **Slide II: 60** : Find every feature element for which the feature is a mountain.
- **Slide II: 61** : Find the name of every mountain.
- **Slide II: 62** : Find the name of every region in which there is a mountain.

The last query was given only in abbreviated form. The full version is more cumbersome:

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

## Abbreviated syntax

The abbreviated syntax is more economical and often (but not always!) more intuitive.

The XPath abbreviations are:

- The syntax **child::** may be omitted from a location step altogether. (The child axis is chosen as default.)
- The syntax **@** is an abbreviation for: **attribute::**
- The syntax **//** is an abbreviation for:  
`/descendant-or-self::node()`
- The syntax **..** is an abbreviation for: **parent::node()**
- The syntax **.** is an abbreviation for: **self::node()**



## Queries and alternatives

Consider again the last query above:

*Find the name of every region in which there is a mountain.*

An alternative location path for this is:

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

Similarly, consider:

*Find the name of countries containing a feature called Everest.*

Two queries for this are:

```
//Feature[text()='Everest']/../../Name/text()
```

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

## One subtle point

A subtle point with XPath is illustrated by the second solution above to:

*Find the name of countries containing a feature called Everest.*

While the given query (repeated below) is correct,

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

the following (natural) attempt would be incorrect:

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

The problem is that the location path `//Feature/text()` starts with a `/` character, and this means that XPath interprets this path as starting at the root node, whereas the path needs to start at the current node.

The omission of a necessary `.` character at the start of a predicate expression is a common source of errors in XPath.

## More on XPath

In practice, when using XPath, one often needs to prefix the location path with a pointer to the given XML document; e.g.,

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

Other features in XPath include: navigation based on document order, position and size of context, treatment of namespaces, a rich language of expressions.

For full details on XPath and XQuery see the W3C specification:

```
http://www.w3.org/TR/xpath
```

A tutorial can be found at:

```
http://www.w3schools.com/xpath/
```

## Part II — Semistructured Data

XML:

**II.1** Semistructured data and XML

**II.2** Structuring XML

**II.3** Navigating XML using XPath

Corpora:

**II.4 Introduction to corpora**

**II.5** Querying a corpus

## Recommended reading

The recommended reading for the material on corpora is:

[CL] Corpus Linguistics  
Tony McEnery & Andrew Wilson  
Edinburgh University Press,  
2nd Edition, 2001

This book is written for a linguistics audience.

Nevertheless, Chapter 2, from the start of chapter to end of §2.2.2, will provide excellent background for the material covered in the lectures.

## Natural language as data

Written or spoken natural language has plenty of *internal structure*: it consists of words, has phrase and sentence structure, etc.

Nevertheless, on a computer, it is represented as a *text file*: simply a sequence of characters.

This is an example of *unstructured data*: the data format itself has no structure imposed on it (other than the sequencing of characters).

Often, however, it is useful to annotate text by marking it up with additional information (e.g. linguistic information, semantic information).

Such marked-up text, is a widespread and very useful form of *semistructured data*.

## What is a corpus?

The word *corpus* (plural *corpora*) is Latin for “body”.

It is used in (both computational and theoretical) linguistics as a word to describe *a body of text*, in particular a body of written or spoken text.

In practice, a *corpus* is a body of written or spoken text, from a particular language variety, that meets the following criteria.

1. sampling and representativeness;
2. finite size;
3. machine-readable form;
4. a standard reference.

## Sampling and representativeness

In linguistics, corpora provide data for *empirical linguistics*

That is, corpora provide data that is used to investigate the nature of linguistic practice (i.e., of real-world language usage), for the chosen language variety

For obvious practical reasons, a corpus can only contain a *sample* of instances of language usage (albeit a potentially large sample)

For such a sample to be useful for linguistic analysis, it must be chosen to be *representative* of the kind of language practice being analysed.

For example, the complete works of Shakespeare would not provide a representative sample for Elizabethan English.



## Finiteness

Furthermore, corpora usually have a fixed *finite* size. It is decided at the outset how the language variety is to be sampled and how much data to include. An appropriate sample of data is then compiled, and the corpus content is fixed.

*N.B. Monitor corpora* (which are beyond the scope of this course) are an exception to the fixed size rule.

While the finite size rule for a corpus is obvious, it contrasts with theoretical linguistics, where languages are studied using *grammars* (e.g. context-free grammars) that potentially generate infinitely many sentences.

## Machine readability

Historically, the word “corpus” was used to refer to a body of printed text.

Nowadays, corpora are almost universally machine (i.e. computer) readable. (Since this is an Informatics course, we are anyway only interested in such corpora.)

Machine-readable corpora have several obvious advantages over other forms:

- They can be huge in size (billions of words)
- They can be efficiently searched
- They can be easily (and sometimes automatically) annotated with additional useful information

## Standard reference

A corpus is often a standard reference for the language variety it represents.

For this, the corpus has to be widely available to researchers.

Having a corpus as a standard reference allows competing theories about the language variety to be compared against each other on the same sample data

The usefulness of a corpus as a standard reference depends upon all the preceding three features of corpora: representativeness, fixed finite size and machine readability.

## Summarizing

In practice, a *corpus* is generally a widely available fixed-sized body of machine-readable text, sampled in order to be maximally representable of the language variety it represents.

Note, however, not every corpus will have all of these characteristics.

## Some prominent English language corpora

- The *Brown Corpus* of American English was compiled at Brown University and published in 1967. It contains around 1,000,000 words.
- The *British National Corpus (BNC)*, published mid 1990's, is a 100,000,000-word text corpus intended to be representative of written and spoken British English from the late 20th century.
- The *American National Corpus (ANC)* is an ongoing project to create an electronic text corpus of written and spoken American English since 1990. The aim is to create a 100,000,000-word corpus. The first release, made available (to subscribers only) in 2003, contains 11,000,000 words and was provided in XML format.
- The *Oxford English Corpus (OEC)* is an English corpus used by the makers of the Oxford English Dictionary. It is the largest text corpus of its kind, containing over 2,000,000,000 words. It is in XML format.

## Two forms of corpus

There are two forms of corpus: *unannotated*, i.e. consisting of just the raw language data, and *annotated*.

Unannotated corpora are examples of *unstructured data*.

Annotated corpora are examples of *semistructured data*.

The four English language corpora on slide II: 77 are all annotated.

Annotations are extremely useful for many purposes. They will play an important role in future lectures.

## Building a corpus

To build a corpus we need to perform two tasks:

- Collect corpus data — this involves *balancing* and *sampling*
- In the case of an annotated corpus, add meta-information — this is called *annotation*

*Balancing* ensures that the linguistic content of a corpus represents the full variety of the language sources that the corpus is intended to provide a reference for. For example, a balanced text corpus includes texts from many different types of source; e.g., books, newspapers, magazines, letters, etc.

*Sampling* ensures that the material is representative of the types of source. For example, sampling from newspaper text: select texts randomly from different newspapers, different issues, different sections of each newspaper.

## Balancing

Things to take into account when balancing:

- *language type*: may wish to include samples from some or all of:
  - edited text (e.g., articles, books, newswire);
  - spontaneous text (e.g., email, Usenet news, letters);
  - spontaneous speech (e.g., conversations, dialogs);
  - scripted speech (e.g., formal speeches).
- *genre*: fine-grained type of material (e.g., 18th century novels, scientific articles, movie reviews, parliamentary debates)
- *domain*: what the material is about (e.g., crime, travel, biology, law);



## Examples of balanced corpora

*Brown Corpus*: a balanced corpus of written American English:

- one of the earliest machine-readable corpora;
- developed by Francis and Kucera at Brown in early 1960's;
- 1M words of American English texts printed in 1961;
- sampled from 15 different genres.

*British National Corpus*: large, balanced corpus of British English.

- one of the main reference corpora for English today;
- 90M words text; 10M words speech;
- text part sampled from newspapers, magazines, books, letters, school and university essays;
- speech recorded from volunteers balanced by age, region, and social class; also meetings, radio shows, phone-ins, etc.

## Comparison of some standard corpora

---

Corpus	Size	Genre	Modality	Language
Brown Corpus	1M	balanced	text	American English
British National Corpus	100M	balanced	text/speech	British English
Penn Treebank	1M	news	text	American English
Broadcast News Corpus	300k	news	speech	7 languages
MapTask Corpus	147k	dialogue	speech	British English
CallHome Corpus	50k	dialogue	speech	6 languages

---

## Pre-processing and annotation

Raw data from a linguistic source can't be exploited directly. We first have to perform:

- *pre-processing*: identify the basic units in the corpus:
  - tokenization;
  - sentence boundary detection;
- *annotation*: add task-specific information:
  - parts of speech;
  - syntactic structure;
  - dialogue structure, prosody, etc.

## Tokenization

*Tokenization*: divide the raw textual data into tokens (words, numbers, punctuation marks).

*Word*: a continuous string of alphanumeric characters delineated by whitespace (space, tab, newline).

*Example*: potentially difficult cases:

- amazon.com, Micro\$oft
- John's, isn't, rock'n'roll
- child-as-required-yuppie-possession  
(As in: “The idea of a child-as-required-yuppie-possession must be motivating them.”)
- cul de sac

## Sentence Boundary Detection

*Sentence boundary detection*: identify the start and end of sentences.

*Sentence*: string of words ending in a full stop, question mark or exclamation mark.

This is correct 90% of the time.

Example: potentially difficult cases:

- Dr. Foster went to Gloucester.
- He said “rubbish!”.
- He lost cash on lastminute.com.

The detection of word and sentence boundaries is particularly difficult for *spoken data*.

## Corpus Annotation

*Annotation*: adds information that is not explicit in the data itself, increases its usefulness (often application-specific).

*Annotation scheme*: basis for annotation, consists of a tag set and annotation guidelines.

*Tag set*: is an inventory of labels for markup.

*Annotation guidelines*: tell annotators (domain experts) how tag set is to be applied; ensure consistency across different annotators.

## Part-of-speech (POS) annotation

*Part-of-speech (POS)* tagging is the most basic kind of linguistic annotation.

Each linguistic token is assigned a code indicating its *part of speech*, i.e., basic grammatical status.

Examples of POS information:

- singular common noun;
- comparative adjective;
- past participle.

POS tagging forms a basic first step in the disambiguation of homographs.

E.g., it distinguishes between the verb “boot” and the noun “boot”.

But it does not distinguish between “boot” meaning “kick” and “boot” as in “boot a computer”, both of which are transitive verbs.

## Example POS tag sets

- CLAWS tag set (used for BNC): 62 tags;
- Brown tag set (used for Brown corpus): 87 tags;
- Penn tag set (used for the Penn Treebank): 45 tags.

Category	Examples	CLAWS	Brown	Penn
Adjective	happy, bad	AJ0	JJ	JJ
Adverb	often, badly	PNI	CD	CD
Determiner	this, each	DT0	DT	DT
Noun	aircraft, data	NN0	NN	NN
Noun singular	woman, book	NN1	NN	NN
Noun plural	women, books	NN2	NN	NN
Noun proper singular	London, Michael	NP0	NP	NNP
Noun proper plural	Australians, Methodists	NP0	NPS	NNPS



## POS Tagging

**Idea:** Automate POS tagging: look up the POS of a word in a dictionary.

**Problem:** POS ambiguity: words can have several possible POS's; e.g.:

Time flies like an arrow. (1)

time: singular noun or a verb;

flies: plural noun or a verb;

like: singular noun, verb, preposition.

**Combinatorial explosion:** (1) can be assigned  $2 \times 2 \times 3 = 12$  different POS sequences.

Need to take sentential context into account to get POS right! A successful approach to this is *probabilistic POS tagging* which can achieve an accuracy of 96–98%.

## Use of markup languages

An important general application of markup languages, such as XML, is to separate *data* from *metadata*.

In a corpus, this serves to keep different types of information apart;

- *Data* is just the raw data.

In a corpus this is the text itself.

- *Metadata* is data about the data.

In a corpus this is the various annotations.

Nowadays, XML is the most widely used markup language for corpora.

The example on the next slide is taken from the BNC XML Edition, which was released only in 2007.

(The previous BNC World Edition was formatted in SGML.)

## Example from the BNC XML Edition

```
<wtext type="FICTION">
  <div level="1">
    <head> <s n="1">
      <w c5="NN1" hw="chapter" pos="SUBST">CHAPTER </w>
      <w c5="CRD" hw="1" pos="ADJ">1</w>
    </s> </head>
    <p> <s n="2">
      <c c5="PUQ"> </c>
      <w c5="CJC" hw="but" pos="CONJ">But</w>
      <c c5="PUN">,</c> <c c5="PUQ"> </c>
      <w c5="VVD" hw="say" pos="VERB">said </w>
      <w c5="NP0" hw="owen" pos="SUBST">Owen</w>
      <c c5="PUN">,</c> <c c5="PUQ"> </c>
      <w c5="AVQ" hw="where" pos="ADV">where </w>
      <w c5="VBZ" hw="be" pos="VERB">is </w>
      <w c5="AT0" hw="the" pos="ART">the </w>
      <w c5="NN1" hw="body" pos="SUBST">body</w>
      <c c5="PUN">?</c> <c c5="PUQ"> </c>
    </s> </p>
    . . . .
  </div>
</wtext>
```

## Aspects of this example

The example is the opening text of J10, a novel by Michael Pearce.

Some aspects of the tagging:

- The **wtext** element stands for *written text*. The attribute **type** indicates the genre.
- The **head** element tags a portion of header text (in this case a chapter heading).
- The **s** element tags sentences. (N.B., a chapter heading counts as a sentence.) Sentences are numbered via the attribute **n**.
- The **w** element tags words. The attribute **pos** is a POS tag, with more detailed POS information given by the **c5** attribute, which contains the CLAWS code. The attribute **hw** represents the *root form* of the word (e.g., the root form of “said” is “say”).
- The **c** element tags punctuation.

## Syntactic annotation (parsing)

*Syntactic annotation*: information about the structure of sentences.

Prerequisite for computing meaning.

Linguists use phrase markers to indicate which parts of a sentence belong together:

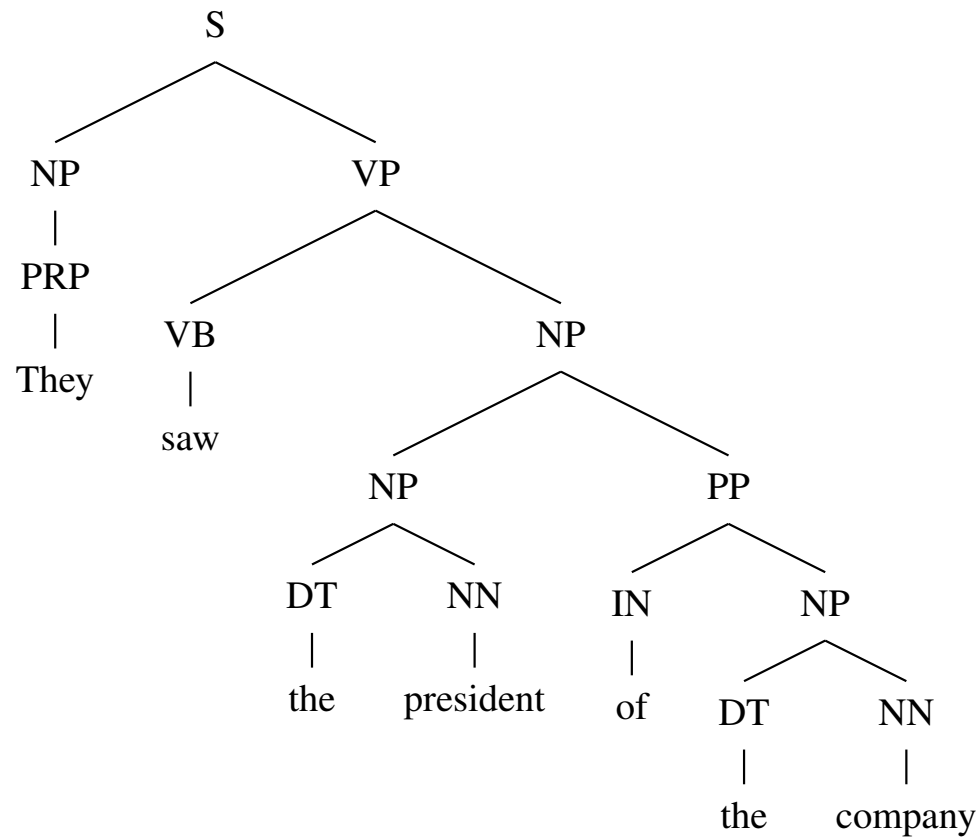
- **noun phrase (NP)**: noun and its adjectives, determiners, etc.
- **verb phrase (VP)**: verb and its objects;
- **prepositional phrase (PP)**: preposition and its NP;
- **sentence (S)**: VP and its subject.

Phrase markers group hierarchically in a *syntax tree*.

Syntactic annotation can be automated. Accuracy: around 90%.

## Example syntax tree

Sentence from the Penn Treebank corpus:



## The same syntax tree in XML:

```
<s>
  <np><w pos="PRP">They</w></np>
  <vp><w pos="VB">saw</w>
    <np>
      <np><w pos="DT">the</w> <w pos="NN">president</w></np>
      <pp><w pos="NN">of</w>
        <np><w pos="DT">the</w> <w pos="NN">company</w></np>
      </pp>
    </np>
  </vp>
</s>
```

Note the conventions used in the above document: **phrase markers** are represented as **elements**; whereas **POS tags** are given as **attribute values**.

**N.B.** The tree on the previous slide is *not* the XML element tree generated by this document.

## Part II — Semistructured Data

### XML:

#### II.1 Semistructured data and XML

#### II.2 Structuring XML

#### II.3 Navigating XML using XPath

### Corpora:

#### II.4 Introduction to corpora

#### **II.5 Querying a corpus**



## Applications of corpora

Answering *empirical questions* in linguistics and cognitive science:

- corpora can be analyzed using statistical tools;
- hypotheses about language processing and language acquisition can be tested;
- new facts about language structure can be discovered.

Engineering *natural-language systems* in AI and computer science:

- corpora represent the data that language processing system have to handle;
- algorithms exist to extract regularities from corpus data;
- text-based or speech-based computer applications can learn automatically from corpus data.

## Extracting data from corpora

To do something useful with corpus data and its annotation, we need to be able to query the corpus to extract the data and information we want.

This lecture introduces:

- The basic notion of a *concordance* in a corpus.
- Statistics are useful for linguistic questions or NLP applications, such as *frequency* and *relative frequency*.
- *Unigrams*, *bigrams* and *n-grams*.
- The linguistic notion of a *collocation*.

## Concordances

*Concordance*: all occurrences of a given word, displayed in context.

More generally, one looks for all occurrences of matches for a given query expression.

- generated by concordance programs based on a user keyword;
- keyword (search query) can specify word, annotation (POS, etc.) or more complex information (e.g., using regular expressions);
- output displayed as keyword in context: matched keyword in the middle of the line, predefined context to left and right.

## Example

A concordance for all forms of the word “*remember*” in a corpus of the complete works of Dickens.

's cellar . Scrooge then <remembered> to have heard that ghost  
, for your own sake , you <remember> what has passed between  
e-quarters more , when he <remembered> , on a sudden , that the  
corroborated everything , <remembered> everything , enjoyed eve  
urned from them , that he <remembered> the Ghost , and became c  
ht be pleasant to them to <remember> upon Christmas Day , who  
its festivities ; and had <remembered> those he cared for at a  
wn that they delighted to <remember> him . It was a great sur  
ke ceased to vibrate , he <remembered> the prediction of old Ja  
as present myself , and I <remember> to have felt quite uncom  
...

## Example

A concordance for all occurrences of “*Holmes*” in a corpus that consists of the Arthur Conan Doyle story *A Case of Identity*.

```
My dear fellow.' ' said Sherlock <Holmes> as we sat on either
a realistic efect,' ' remarked <Holmes>.  ``This is wanting in the
said <Holmes>, taking the paper and glancing his eye down
``I have seen those symptoms before,' ' said <Holmes>, throwing
merchant-man behind a tiny pilot boat.  Sherlock <Holmes> welcomed
You've heard about me, Mr. <Holmes>,' ' she cried, ``else how
...
```

## Frequencies

Frequency information obtained from corpora is often useful for answering scientific or engineering questions.

*Token count  $N$* : number of tokens (words, punctuation marks, etc.) in a corpus (i.e., size of the corpus).

*Type count*: number of *different* tokens in a corpus.

*Absolute frequency  $f(t)$  of a type  $t$* : number of tokens of type  $t$  in a corpus.

*Relative frequency of a type  $t$* : absolute frequency of  $t$  normalized by the token count, i.e.,  $f(t)/N$ .

## Frequencies (example)

The British National Corpus (BNC) is an important reference.

Let's compare some counts from the BNC with counts from our sample corpus *A Case of Identity*

	BNC	A Case of Identity
Token count $N$	100,000,000	7,006
Type count	636,397	1,621
$f(\text{Holmes})$	890	46
$f(\text{Sherlock})$	209	7
$f(\text{Holmes})/N$	.0000089	.0066
$f(\text{Sherlock})/N$	.00000209	.000999

## Unigrams

We can now ask questions such as: what are the most frequent words in a corpus?

- Count absolute frequencies of all word types in the corpus;
- tabulate them in an ordered list;
- results: list of *unigram* frequencies (frequencies of individual words).

The next slide compares unigram frequencies for BNC and *A Case of Identity*.



## Unigrams (example)

BNC		A Case of Identity	
6,184,914	the	350	the
3,997,762	be	212	and
2,941,372	of	189	to
2,125,397	a	167	of
1,812,161	in	163	a
1,372,253	have	158	I
1,088,577	it	132	that
917,292	to	117	it

**N.B.** The article “the” is the most frequent word in both corpora; prepositions like “of” and “to” appear in both lists; etc.

## $n$ -grams

The notion of unigram can be generalized:

- *bigrams* — pairs of adjacent words
- *trigrams* — triples of adjacent words
- *$n$ -grams* —  $n$ -tuples of adjacent words.

As the value of  $n$  increases, the units become more linguistically meaningful.

*n*-grams (example)

Compute the most frequent *n*-grams in *A Case of Identity*, for  $n = 2, 3, 4$ .

bigrams		trigrams		4-grams	
40	of the	5	there was no	2	very morning of the
23	in the	5	Mr. Hosmer Angel	2	use of the money
21	to the	4	to say that	2	the very morning of
21	that I	4	that it was	2	the use of the
20	at the	4	that it is	2	the King of Bohemia

**N.B.** *n*-gram frequencies get smaller with increasing *n*. As more word combinations become possible, there is increased *data sparseness*.

## Example

A concordance for all occurrences of bigrams in the Dickens corpus in which the second word is “*tea*” and the first is an adjective.

This query exploits the POS tagging of the corpus to search for adjectives.

```
now , notwithstanding the <hot tea> they had given me before  
' ' Shall I put a little <more tea> in the pot afore I go ,  
o moisten a box-full with <cold tea> , stir it up on a piece  
tween eating , drinking , <hot tea> , devilled grill , muffi  
e , handed round a little <stronger tea> . The harp was there ; t  
e so repentant over their <early tea> , at home , that by eigh  
rs. Sparsit took a little <more tea> ; and , as she bent her  
s illness ! Dry toast and <warm tea> offered him every night  
of robing , after which , <strong tea> and brandy were administ  
rsty . You may give him a <little tea> , ma'am , and some dry t
```

## Collocations

*Collocation*: a sequence of words that occurs ‘atypically often’ in language usage

Examples:

- *run amok*: the verb “run” can occur on its own, but “amok” can’t.
- *strong tea*: sounds much better than “powerful tea” although the literal meanings are much the same.
- Phrasal verbs such as *make up* or *make off* or *make out* (but not, for example, “make in”).
- *rancid butter*, *bitter sweet*, *over and above*, etc.

**N.B.** The inverted commas around ‘atypically often’ are because we shall eventually need statistical ideas to make this precise.

## Identifying collocations

**Task:** automatically identify collocations in a large corpus.

For example collocations with the word *tea* (see III: 108).

- *strong tea* occurs in the corpus.

This is a collocation.

- *powerful tea*, in fact, does not.

- However, *more tea* and *little tea* also occur in the corpus.

These are not collocations. These word sequences do not occur with an *atypically* common frequency.

**Problem:** How do we detect when a bigram (or  $n$ -gram) is a collocation?

## Looking at the data

The next slide lists the frequencies of the most common bigrams, in the Dickens Corpus, in which the first word is “*strong*”.

For comparison, the frequencies of the most common bigrams in which the first word is “*powerful*” are also given.

strong	and	31	powerful	effect	3
	enough	16		sight	3
	in	15		enough	3
	man	14		mind	3
	emphasis	11		for	3
	desire	10		and	3
	upon	10		with	3
	interest	8		enchanter	2
	a	8		displeasure	2
	as	8		motives	2
	inclination	7		impulse	2
	tide	7		struggle	2
	beer	7		grasp	2



## Filtering collocations

The bigram table shows:

- Neither *strong tea* nor *powerful tea* are frequent enough to make it into the top 13.
- Potential collocations for *strong*: e.g., *strong desire*, *strong inclination*, and *strong beer*;
- Potential collocations for *powerful*: e.g., *powerful effect*, *powerful motives*, and *powerful struggle*;
- Problem: The bigrams *strong and*, *strong enough*, *powerful for*, are highly frequent. These are not collocations.
- To distinguish collocations from non-collocations, we need to filter out ‘noise’.

## The need for statistics

**Problem:** Words like *for* and *and* are highly frequent on their own: they occur with *tea* by chance.

**Solution:** use statistical testing to detect when the frequency of a bigram is *atypically high* given the frequencies of its constituent words.

In general, statistical tools offer powerful methods for the analysis of all types of data. In particular, they provide the principal approach to the quantitative (and qualitative) analysis of *unstructured data*.

We shall return to the problem of finding collocations in Part III of the course, when we have appropriate statistical tools at our disposal.

## Searching for concordances

The concordances in this lecture were produced using a dedicated program for searching for concordances, the *Corpus Query Processor (CQP)*.

CQP is query engine which searches corpora based on user queries over words, parts of speech, or other markup.

It uses *regular expressions* to formulate queries. This makes the CQP query language very powerful (N.B. This is the second time we have found an application for regular expressions in Data & Analysis.)

An alternative to using a dedicated concordance program is to use XML query technology (XPath and XQuery) to search any corpus implemented in XML.

## Corpora in Informatics

Corpora are used extensively in two areas of informatics:

- *Natural Language Processing (NLP)* builds computer systems that understand or produce text. Example applications that rely on corpus data include:
  - *Summarization*: take a text and compress it, i.e., produce an abstract or summary. Example: Newsblaster.
  - *Machine Translation (MT)*: take a text in a source language and turn it into a text in the target language. Example: Babel Fish.
- *speech processing* develops systems that understand or produce spoken language.

The techniques applied rely on probability theory, information theory and machine learning to extract statistical regularities from corpora.

Example translation by AltaVista Babel Fish.

*O, my love is like a red, red rose,  
That is newly sprung in June.*

Robert Burns (1759–1796)

English → Italian:

*La O, il mio amore è come un rosso, colore rosso è aumentato,  
che recentemente è balzato in giugno.*

Italian → English:

*Or, my love is like a red one, red color is increased,  
than recently it is jumped in June.*

Fortunately, Machine-translation research has made progress since Babel Fish.