

Report Project 1 : XML Schema Definition

Rémy Detobel & Nathan Liccardo

March 29, 2018

1 Introduction

The aims of this document will be to describe and explain all the choices that we made for our project. To call back, the main goal of this assignment was to create an XSD file which contains a specific XML Schema Definition. The first part of this report will be dedicated to the variables that we defined. In a second step we will show, by means of trees, how the requested structure has been implemented.

2 Variables

Many variables are strings whose formatting is quite free. This is the case, for :

- title
- publisher
- abstract
- edition
- author
- editor

Most of the other fields are integers, such as

- volume
- number
- price
- impact

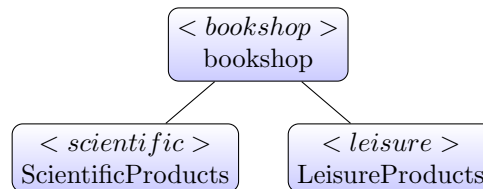
Three types remain to be defined : the year, the ISBN and the genre. For the first one we defined a global type (named “yearType” which will be reused later), as well as a classic element type “year”. For ISBN, many formats are possible, more information is given in the section 2.1. For genre, we have limited the possible entries to the following list : thriller, horror, sci/fi, romance, literature.

The genre and the ISBN are based on the string type that we have restrict, for the first via a list listing all the possibilities and for the second with a pattern (as described in the point 2.1).

2.1 ISBN

3 Structure of the schema

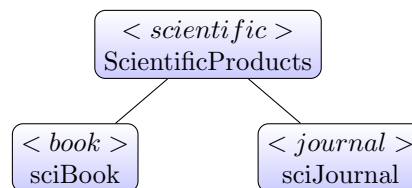
For this project, we were assigned to write an XSD file which define a book shop. This shop is separated into two parts : scientific products and leisure products. The following tree represent this first relation :



In this tree, each node contains the XML tag but also the type of the element. So here we can see that the “bookshop” type does have two children : `scientific` and `leisure`. These two types will be detailed in the following points.

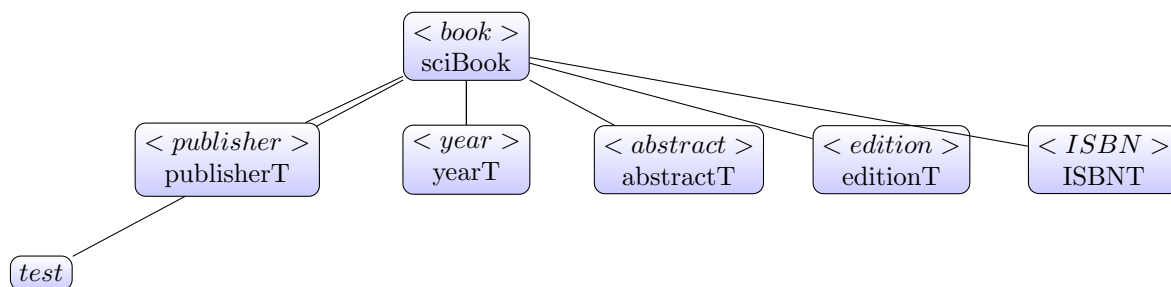
3.1 Scientific products

Scientific products are separated into two sub-categories. The first one define scientific books and the second one scientific journals. Those links can be represented as follow :



These two new type (`sciBook` and `sciJournal`) are a little more complex than what we have seen so far.

3.1.1 Scientific book



3.1.2 Scientific journal

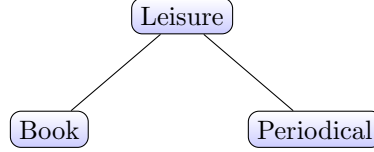
Based on this, we can define multiple relations. Those relations will be used to represent complex types in the XSD file. As shown in the previous tree, we can define three complex types named `scientificType`, `bookType` and `journalType` :

- `bookshopType` \rightarrow (`scientific`, `scientificType`)
- `scientificType` \rightarrow (`book`, `booktype1`), (`journal`, `journalType`)

For the moment, we are not going to define the `bookshopType`. In fact, we do not have enough informations to define it correctly now. We are so assuming that it is an abstract type. Regarding the two final complex types, called `bookType` and `journalType`, a complete explanation of these will be given in the second part dedicated to hypotheses and choices.

3.2 Leisure products

The second section of a book shop is dedicated to leisure products. As for scientific one, products are separated into two sub-categories. The first one is used to define leisure books and the second one leisure periodicals. Once again , those links can be illustrated by means of a tree :



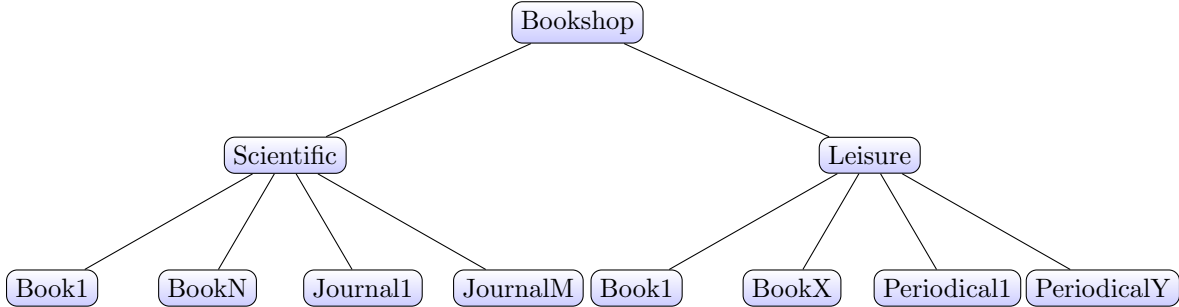
Following this tree, we can obtain two relations. These last are used to define an abstract bookshop (bookshopType) and a final leisure (leisureType). Finally, we obtain the next definitions :

- $\text{bookshopType} \rightarrow (\text{leisure}, \text{leisureType})$
- $\text{leisureType} \rightarrow (\text{book}, \text{booktype2}), (\text{periodical}, \text{periodicalType})$

As said in the previous sub-section, bookshopType is actually not completely defined. For this reason, there exist some incompatibilities between the scientific and leisure part. Notice that leisure books and scientific books are not using same types. In fact we will see in the next section that they have two different definitions.

3.3 Complete schema

The two previous parts were focused on the scientific and leisure definitions. We have now to merge those definitions to obtain the final tree. In fact, a bookshop is a combination of scientific and leisure products. Notice that for each element (book, journal, ...) we can also have multiple instances of it. Finally, we obtain the following tree :



Previously, we defined the bookshopType as an abstract one. In fact, this type was a sequence of scientificType and leisureType and could not be defined in one time. As we changed the definition of a bookshopType but also of a scientificType and a leisureType, we obtain three new relations :

- $\text{bookshopType} \rightarrow (\text{scientific}, \text{scientificType}), (\text{leisure}, \text{leisureType})$
- $\text{scientificType} \rightarrow (\text{book}, \text{booktype1})^*, (\text{journal}, \text{journalType})^*$
- $\text{leisureType} \rightarrow (\text{book}, \text{booktype2})^*, (\text{periodical}, \text{periodicalType})^*$

For the moment, we only defined links between elements. All those are created by defining complex types with some elements inside. Each of these elements are themselves referring to other complex type excepted for final types (Book, Journal, ...). The second part of this report will be focused on explaining and illustrating all the choices that we made during the implementation. For this reason, we will also detail the structure of each final types.

4 Hypotheses and choices

All along the previous section, we only described the global structure of the schema. We will now focus on what we called previously "final types". There exist four different kinds of this specific element : bookType1, journalType, bookType2 and periodicalType. As said in the previous section, bookType1 and bookType2 are both identified by the <book> tag but are not referring to the same type. This can be achieved due to the scope definition. In fact, bookType1 is defined inside the scientific scope unlike bookType2 which is defined inside the leisure scope. Regarding to the two last types, journalType is identified by the <journal> tag (scientific scope) and periodicalType by the <periodical> tag (leisure scope). This section will be separated into two sub-sections : the XSD's structure and the final types structures.

4.1 XSD's structure

Our file has been split into five distinct parts. Each part is used to represent one level of the XML Schema. Indeed, the order of declarations is made from bottom to top level. Here is the used structure :

1. Simple types : This first section is used to define each atomic elements. This is mainly useful to reduce duplicate definitions and code. For example, we define the title element, the publisher element, etc.
2. Scientific part : Define all complex types which belong to the scientific product tree. We can see that scientificBook and scientificJournal are defined at the end of this section. Detail about this part will be given in the next section.
3. Leisure part : As the scientific part, leisure part is used to define all complex types which belong to the leisure tree. Once again, details about this part will be given later.
4. Products part : Define two complex types, scientific product and leisure product. Both can be seen as root elements from the two previous parts.
5. Bookshop part : Is used to merge both products elements. This section also define the root element of the xml file named <b:bookshop>.

4.2 Final types

This second section will be focused on the structure of each final types. As said before, there exist four different final types : bookType1, journalType, bookType2 and periodicalType.

4.2.1 BookType1

According to the assignment, a scientific book must be composed of : a title, a list of authors xor a list of editors, a publisher and a year of publication. Optionally, a book can have : an abstract, an edition and an ISBN. Formally, we obtain the following definition :

- $\text{bookType1} \rightarrow (\text{title}, \text{titleT}), (\text{authors}, \text{authorsT})^+ \mid (\text{editors}, \text{editorsT})^+, (\text{publisher}, \text{publisherT}), (\text{year}, \text{yearT}), (\text{abstract}, \text{abstractT})?, (\text{edition}, \text{editionT})?, (\text{ISBN}, \text{ISBNT})?$

All the types, excepted authorsT and editorsT, used in the previous section are defined in the simple types section. Regarding to authorsT and editorsT, their definitions are more complex. Indeed, they are themselves complex types definitions. We made the arbitrary choice that a list of authors (or editors) must be structured as follow :

```
<authors>
  <author> ... </author>
  ....
  <author> ... </author>
</authors>
```

Which implies that author and editor tags are defined by authorT and editorT. These new types are finally defined in the simple type section (as string elements).

4.2.2 JournalType

A scientific journal has also two parts : a mandatory one and an optional one. For the mandatory part, we must have : a title, a volume, a number, a list of authors xor a list of editors and a year. Concerning to the optional part, a journal type can have a publisher and an impact. In addition, a journal must also have a list of articles. Finally, we obtain this relation :

- $\text{journalType} \rightarrow (\text{title}, \text{titleT}), (\text{volume}, \text{volumeT}), (\text{number}, \text{numberT}), (\text{authors}, \text{authorsT})^+ \mid (\text{editors}, \text{editorsT})^+, (\text{year}, \text{yearT}), (\text{publisher}, \text{publisherT})?, (\text{impact}, \text{impactT})?, (\text{articles}, \text{articlesT})$

Were we have two new complex types : impact and article. All the other one are defined in the simple type section excepted editors and authors which were defined previously. An impact is defined with two elements namely a year and an impact factor. For this reason, we decided to structure this complex type as follow :

```
<impact year="XXXX"> ... </impact>
```

Notice that another solution was to create a complex type (impact) composed by a year tag and a factor tag. Regarding to the articles type, we had to define a list of articles. To do that, we defined a complex type composed itself by a list of article tags. Here is the XML logic :

```
<articles>
  <article> ... </article>
  ...
  <article> ... </article>
</articles>
```

Inside each article, we must find those elements : a title, an author, an interval of pages xor a number corresponding to the article. This can be formalised by the following relation :

- $\text{articlesT} \rightarrow (\text{article}, \text{articleT})^+$
- $\text{articleT} \rightarrow (\text{title}, \text{titleT}), (\text{author}, \text{authorT}), (\text{pages}, \text{pagesT}) \mid (\text{number}, \text{numberT})$

Note that this type contains only one author but also that we use the complex type pagesT and the simple type numberT (define in the first section). PagesT is a very simple complex type composed of two integers corresponding to the start and end page.

4.2.3 BookType2

The second book type is defined inside the leisure scope. For this reason, the composition of this new book type will differ from bookType1. Inside this new book type we must find : a title, a list of authors, a publisher, a year and a genre. Optionally, we can add an edition and a number. Formally, we obtain the following relation :

- $\text{bookType2} \rightarrow (\text{title}, \text{titleT}), (\text{authors}, \text{authorsT}), (\text{publisher}, \text{publisherT}), (\text{genre}, \text{genreT}), (\text{edition}, \text{editionT})?, (\text{number}, \text{numberT})?$

We were asked to restrict the value of a genre. To do that we used a restriction tag with multiple enumerations inside. All the other types were defined previously or are defined in the simple type section.

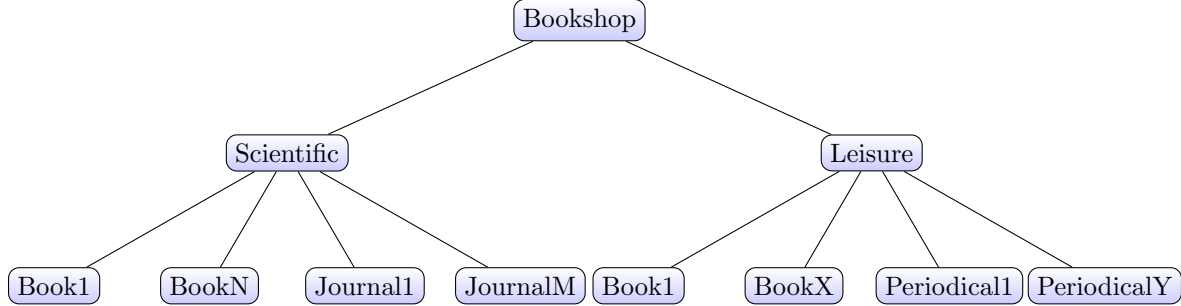
4.2.4 PeriodicalType

This last complex type is only composed by a sequence of three simple type elements. As these elements are very simple, we will not give any more informations about them. Here is the formal definition of a periodicalType :

- $\text{periodicalType} \rightarrow (\text{title}, \text{titleT}), (\text{price}, \text{priceT}), (\text{publisher}, \text{publisherT})$

4.3 Complete schema and relations

To complete this report, we can now show all the relations corresponding to the final tree. Note that those relations are now correctly named (as same as in the XSD file) and that we order relations from top to bottom (reverse file order). Due to their trivial aspects, we didn't add simple type (texts, numbers, ...)



- bookshopType \rightarrow (scientific, scientificType),(leisure, leisureType)
- scientificType \rightarrow (book, booktype1)*, (journal, journalType)*
- leisureType \rightarrow (book, booktype2)*, (periodical, periodicalType)*
- bookType1 \rightarrow (title, titleT), (authors, authorsT)⁺ | (editors, editorsT)⁺, (publisher, publisherT), (year, yearT), (abstract, abstractT)?, (edition, editionT)?, (ISBN, ISBNNT)?
- journalType \rightarrow (title, titleT), (volume, volumeT), (number, numberT), (authors, authorsT)⁺ | (editors, editorsT)⁺, (year, yearT), (publisher, publisherT)?, (impact, impactT)?, (articles, articlesT)
- articlesT \rightarrow (article, articleT)⁺
- articleT \rightarrow (title, titleT), (author, authorT), (pages, pagesT) | (number, numberT)
- bookType2 \rightarrow (title, titleT), (authors, authorsT), (publisher, publisherT), (genre, genreT), (edition, editionT)?, (number, numberT)?
- periodicalType \rightarrow (title, titleT), (price, priceT), (publisher, publisherT)