

Big Data – Exercises

Fall 2021 – Week 6 – ETH Zurich

Data Models

Reading:

- (Mandatory) Melnik, S., et al. (2011). Dremel: Interactive Analysis of Web-Scale Datasets. In CACM, 54(6). [[DOI](#)]
- (Recommended) M. Droettboom, Understanding JSON Schema [[online](#)]
- (Recommended) Harold, E. R., & Means, W. S. (2004). XML in a Nutshell. [Available in the ETH library] [[online](#)] (Chapter 17 on XML Schema, except 17.3 on namespaces)
- (Optional) White, T. (2015). Hadoop: The Definitive Guide (4th ed.). O'Reilly Media, Inc. [Available in the ETH library] [[online](#)] (Chapters 12 (Avro) and 13 (Parquet))

This exercise will consist of three main parts:

- XML data models
- XML schemas
- JSON Schemas
- JSound
- Dremel

1. XML Data Models – Information Sets

XML "Information Set" provides an abstract representation of an XML document—it can be thought of as a set of rules on how one would draw an XML document on a whiteboard.

An XML document has an information set if it is well-formed and satisfies the namespace constraints. There is no requirement for an XML document to be valid in order to have an information set. An information set can contain up to eleven different types of information items, e.g., the document information item (always present), element information items, attribute information item, etc.

Task 1.1

Draw the Information Set trees for the following XML documents. You can confine your trees to only have the following types of information items: *document information item*, *elements*, *character information items*, and *attributes*.

Document 1

```
1 <Burger>
2   <Bun>
3     <Pickles/>
4     <Cheese origin="Switzerland" />
5     <Patty/>
6   </Bun>
7 </Burger>
```

Document 2

```
1 <catalog>
2   <!-- A list of books -->
3   <author>Gambardella, Matthew</author>
4   <title>XML Developer's Guide</title>
5   <genre>Computer</genre>
6   <price>44.95</price>
7   <publish_date version='hard' version2='soft'>2000-10-01</publish_date>
8 </book>
9 </catalog>
```

Document 3

```
1 <eth date="11.11.2006">
2   <date>16.11.2017</date>
3   <president since="2020">Prof. Dr. Joël Mesot</president>
4   <rector>Prof. Dr. Sarah M. Springman</rector>
5 </eth>
```

2. XML Schemas

In this task we will explore XML Schemas in detail. An XML Schema describes the structure of an XML document.

The purpose of an XML Schema is to define the legal building blocks of an XML document:

- the elements and attributes that can appear in a document
- the number of (and order of) child elements
- data types for elements and attributes
- default and fixed values for elements and attributes

When you open an XML Schema in oXygen, you can switch to its graphical representation, by choosing the "Design" mode at the bottom of the document pane; "Text" mode shows the XML Schema as an XML document.

Task 2.1

Match the following XML documents to XML Schemas that will validate them. Match them manually then validate with oXygen.

Document 1

```
1 <happiness xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
2   xsi:noNamespaceSchemaLocation="Schema.xsd"/>
```

Document 2

```
1 <happiness xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
2   xsi:noNamespaceSchemaLocation="Schema.xsd">
3   <health/>
4   <friends/>
5   <family/>
6 </happiness>
```

Document 3

```
1 <happiness xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
2   xsi:noNamespaceSchemaLocation="Schema.xsd">
3   3.141562
4 </happiness>
```

Document 4

```
1 <happiness xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
2   xsi:noNamespaceSchemaLocation="Schema.xsd">
3   <health value="100"/>
4   <friends/>
5   <family/>
6 </happiness>
```

Document 5

```
1 <happiness xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"  
2   xsi:noNamespaceSchemaLocation="Schema.xsd">  
3   <health/>  
4   <friends/>  
5   <family/>  
6   But perhaps everybody defines it differently...  
7 </happiness>
```

Schema 1

```
1 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
2   <xs:element name="happiness">
3     <xs:complexType>
4       <xs:sequence>
5         <xs:element name="health"/>
6         <xs:element name="friends"/>
7         <xs:element name="family"/>
8       </xs:sequence>
9     </xs:complexType>
10  </xs:element>
11 </xs:schema>
```

Schema 2

```
1 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
2   <xs:element name="happiness">
3     <xs:complexType mixed="true">
4       <xs:sequence>
5         <xs:element name="health"/>
6         <xs:element name="friends"/>
7         <xs:element name="family"/>
8       </xs:sequence>
9     </xs:complexType>
10  </xs:element>
11 </xs:schema>
```

Schema 3

```
1 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
2   <xs:element name="happiness" type="xs:decimal"/>
3 </xs:schema>
```

Schema 4

```
1 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
2   <xs:element name="happiness">
3     <xs:complexType>
4       <xs:sequence/>
5     </xs:complexType>
6   </xs:element>
7 </xs:schema>
```

Schema 5

```
1 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
2   <xs:element name="happiness">
3     <xs:complexType>
4       <xs:sequence>
5         <xs:element name="health">
6           <xs:complexType>
7             <xs:attribute name="value" type="xs:integer"
8             use="required"/>
9           </xs:complexType>
10          </xs:element>
11          <xs:element name="friends"/>
12          <xs:element name="family"/>
13        </xs:sequence>
14      </xs:complexType>
15    </xs:element>
16  </xs:schema>
```


Task 2.2

The [Great Language Game](#) is a game in which you are given a voice clip to listen, and you are asked to identify the language in which the person was speaking. It is a multiple-choice question—you make your choice out of several alternatives.

The following XML document presents a user's attempt at answering a single question in the game: it contains the identifier of the voice clip, the choices presented to the player, and the player's response.

Provide an XML Schema which will validate this document:

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <attempt country="AU" date="2013-08-19">
3   <voiceClip>48f9c924e0d98c959d8a6f1862b3ce9a</voiceClip>
4   <choices>
5     <choice>Maori</choice>
6     <choice>Mandarin</choice>
7     <choice>Norwegian</choice>
8     <choice>Tongan</choice>
9   </choices>
10  <target>Norwegian</target>
11  <guess>Norwegian</guess>
12 </attempt>
13
```

Task 2.3

Continuing the topic of the Great Language Game, provide an XML Schema which will validate the following document:

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <attempts>
3   <attempt country="AU" date="2013-08-19">
4     <voiceClip>48f9c924e0d98c959d8a6f1862b3ce9a</voiceClip>
5     <choices>
6       <choice>Maori</choice>
7       <choice>Mandarin</choice>
8       <choice>Norwegian</choice>
9       <choice>Tongan</choice>
10    </choices>
11    <target>Norwegian</target>
12    <guess>Norwegian</guess>
13  </attempt>
14  <attempt country="US" date="2014-03-01">
15    <voiceClip>5000be64c8cc8f61dda50fca8d77d307</voiceClip>
16    <choices>
17      <choice>Finnish</choice>
18      <choice>Mandarin</choice>
19      <choice>Scottish Gaelic</choice>
20      <choice>Slovak</choice>
21      <choice>Swedish</choice>
22      <choice>Thai</choice>
23    </choices>
24    <target>Slovak</target>
25    <guess>Slovak</guess>
26  </attempt>
27  <attempt country="US" date="2014-03-01">
28    <voiceClip>923c0d6c9e593966e1b6354cc0d794de</voiceClip>
29    <choices>
30      <choice>Hungarian</choice>
31      <choice>Sinhalese</choice>
32      <choice>Swahili</choice>
33    </choices>
34    <target>Hungarian</target>
35    <guess>Sinhalese</guess>
36  </attempt>
37 </attempts>
```

Task 2.4

Let us now solve the reverse problem. Given the following XML Schema, provide a valid instance document:

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
3   <xs:element name="movies">
4     <xs:complexType>
5       <xs:sequence maxOccurs="unbounded" minOccurs="0">
6         <xs:element name="Movie">
7           <xs:complexType>
8             <xs:sequence>
9               <xs:element name="title" type="xs:string"/>
10              <xs:element name="year" type="xs:gYear"/>
11              <xs:element name="_director">
12                <xs:complexType>
13                  <xs:sequence/>
14                  <xs:attribute name="name" type="xs:string"/>
15                </xs:complexType>
16              </xs:element>
17              <xs:choice minOccurs="1" maxOccurs="unbounded">
18                <xs:element name="comment">
19                  <xs:complexType>
20                    <xs:simpleContent>
21                      <xs:extension base="xs:string">
22                        <xs:attribute name="lang"
23                          type="xs:string"/>
24                      </xs:extension>
25                    </xs:simpleContent>
26                  </xs:complexType>
27                <xs:element name="newcomment" type="xs:string"/>
28              </xs:choice>
29            </xs:sequence>
30            <xs:attribute name="id" type="xs:ID"/>
31          </xs:complexType>
32        </xs:element>
33      </xs:sequence>
34    </xs:complexType>
35  </xs:element>
36 </xs:schema>
```

3. JSON Schemas

JSON Schema is a vocabulary that allows you to annotate and validate JSON documents. It is used to:

- Describe your existing data format(s).
- Provide clear human- and machine- readable documentation.
- Validate data, i.e., automated testing, ensuring quality of client submitted data.

Task 3.1

Provide an JSON Schema which will validate the following document.

```
1 {  
2   "firstName": "John",  
3   "lastName": "Doe",  
4   "age": 21  
5 }
```

Provide an JSON Schema which will validate the following document.

The JSON Schema has to check for the following properties:

- The price of a product has to be strictly positive.
- Tags are describing the product and necessary for a proper product description. We need at least one tag per product and each tag should be unique.
- The "productId", "productName" and the "price" should always be contained in a valid JSON document.

```
1  {
2    "productId": 1,
3    "productName": "An ice sculpture",
4    "price": 12.50,
5    "tags": [ "cold", "ice" ],
6    "dimensions": {
7      "length": 7.0,
8      "width": 12.0,
9      "height": 9.5
10   }
11 }
```

Task 3.3

Based on the given Json schema, can you give an instance of it?

HINT: We defined array of things.

```
1  {
2    "$id": "https://example.com/arrays.schema.json",
3    "$schema": "http://json-schema.org/draft-07/schema#",
4    "description": "A representation of a person, company, organization, or place",
5    "type": "object",
6    "properties": {
7      "fruits": {
8        "type": "array",
9        "items": {
10         "type": "string"
11       }
12     },
13     "vegetables": {
14       "type": "array",
15       "items": { "$ref": "#/definitions/veggie" }
16     }
17   },
18   "definitions": {
19     "veggie": {
20       "type": "object",
21       "required": [ "veggieName", "veggieLike" ],
22       "properties": {
23         "veggieName": {
24           "type": "string",
25           "description": "The name of the vegetable."
26         },
27         "veggieLike": {
28           "type": "boolean",
29           "description": "Do I like this vegetable?"
30         }
31       }
32     }
33   }
34 }
```

4. JSound

JSound is a vocabulary that allows you to validate JSON documents. It employs a very simple and intuitive JSON-like syntax.

Task 4.1

Repeat the exercise in 3.1 but now produce a JSound schema which will validate the following document.

```
1 {  
2   "firstName": "John",  
3   "lastName": "Doe",  
4   "age": 21  
5 }
```

Task 4.2

Build a valid JSON document based on the following JSound schema.

```
1  {
2    "id": "integer",
3    "who": [{
4      "name": "string",
5      "type": "string",
6      "preferred": "boolean"
7    }],
8    "year_of_birth": "integer",
9    "living": "boolean"
10 }
```


5. Document Validation

XML, JSON, and other document formats have many different tools for validation.

Associate each document type with the type of schema that can be used to validate it:

1	Document Type:	Schema Type:
2		
3	A. XML	1. XML Schema
4	B. JSON	2. DTD
5	C. Protocol buffers	3. Schematron
6	D. XHTML	4. RelaxNG
7		5. JSON Schema
8		6. JSound
9		7. Kwalify
10		8. The XML Schema for XHTML
11	documents	9. proto3
12		10. No schema at all

6. Dremel

Dremel is a query system developed at Google for deriving data stored in a nested data format such as XML, JSON, or Google Protocol Buffers into column storage, where it can be analyzed faster.

For this exercise, you will apply the algorithm described in the [Dremel paper](#) to convert from a nested data format with an associated schema into column storage.

You are given the following Google Protocol Buffer schema:

```
1  message Catalog {
2      repeated group Book {
3          required string Title;
4          repeated string Author;
5          optional uint32 Year;
6          repeated group Version {
7              required string Type;
8              optional uint32 Pages;
9          }
10     }
11 }
```

Use this to convert the following message to column storage:

```
1  Book
2      Title: "Because Internet"
3      Author: "McCulloch, Gretchen"
4      Year: 2019
5      Version
6          Type: "soft"
7  Book
8      Title: "Hitchhiker's Guide to the Galaxy"
9      Author: "Adams, Douglas"
10     Year: 1979
11  Book
12     Title: "XML in a Nutshell"
13     Author: "Harold, Elliotte Rusty"
14     Author: "Means, Scott W."
15     Year: 2005
16     Version
17         Type: "soft"
18         Pages: 718
```

Fill in each of the following tables storing the data for each leaf field, where each row contains the value, the repetition levels, and the definitions levels for the given entry.

HINTS: the definitions of repetition levels and definitions levels are presented in the section 4.1 in the [Dremel paper](#)

Catalog.Book.Title

Catalog.Book.Author

Catalog.Book.Year

Catalog.Book.Version.Type

Catalog.Book.Version.Pages