

Exercise - Data Exchange Formats

Web Data Integration IE683
University of Mannheim, Germany

September 2023

In the following three exercises, you are asked to write Java code for reading data from XML, JSON and RDF files and for querying the data using the XPath and SPARQL query languages. Each subsection is dedicated to one of the three data exchange formats. The tasks are rather basic and the goal is to refresh your knowledge in Java in general and in particular in parsing those formats.

1 How to start with Java?

After installing Java 11, download and unzip the DataParser project which can be found on the course web page. Import the project as a Maven project in Eclipse. In order to do so, select inside Eclipse File/Import/Import Existing Maven Project and point Eclipse to the DataParser directory containing the Maven file pom.xml. Now inspect the structure of the project. Inside the src.main.java folder, you can find three packages. Each package is dedicated to one data format: JSON, RDF and XML. For each format, we have prepared some Java classes which will help you solve the tasks of this exercise.

2 XML

This subsection is dedicated to the XML format. In particular, you are asked to perform XPath queries on the *Mondial* dataset¹. This dataset includes world geographic information integrated from the *CIA World Factbook*, the *International Atlas* and the *TERRA database*, to name just the pre-dominant sources. Please inspect the document manually (using a text editor) in order to explore the structure. You can also have a look at the *w3school XPath tutorial*² to solve the following tasks.

¹The file can be downloaded from ILIAS but is also available here: <http://aiweb.cs.washington.edu/research/projects/xmltk/xmldata/data/mondial/mondial-3.0.xml>

²https://www.w3schools.com/xml/xpath_intro.asp

2.1 Starting point

Locate the `XMLReader.java` class in the package `de.dwslab.lecture.wdi.xml`. This class contains a main method which parses the `mondial-3.0.xml` file using the *JAXP* library and selects its root node.

2.2 Mondial — Print XML node information

TASK 1: In order to get started with parsing of XML files in Java, please write a Java class which reads the `mondial-3.0.xml` file using the JAXP library and prints the name of the root element of the XML document. In order to make your start easier, the `DataParser` project contains a file `WDI_20200_RG2_Template1.java` which already includes all the imports for solving the task, as well as comments explaining to you which steps need to be performed in which order.

2.3 Mondial — Schema Inspection

Now that we have written our parser and explored the root node we can start digging deeper into the XML file.

TASK 2: Adapt the class from the previous task so that a unique list of all nodes below the root node is printed. (Hint: In order to prevent that the same node name is printed multiple times when iterating over the `NodeList`, you can use a Java `HashSet` for remembering which node names have already been printed.)

2.4 Mondial — Basic XPath

Now that we got an idea about the structure of the XML we are interested in the content.

TASK 3: Adapt the solution of the previous task in the way that it prints the names of all countries which belong to the continent with the name **Europe**. (Hint: Have a look at the schema of the node `country` to see how it is linked to the continent.)

2.5 Mondial — XPath Predicates I

With the solution of the previous task, we are now able to get the countries for a selected continent. In the next step, we want to extend this query so that we can get countries which belong to two continents.

TASK 4: Extend the XPath for the former task in order to retrieve only countries which are part of **Europe** and **Asia**.

2.6 Mondial — XPath Predicates II

In a final step, we want to gather all attributes from a selection of nodes, without explicitly knowing their names.

TASK 5: Extend the solution of the former task in order to navigate (using XPath) to the `country` node and print all attribute names and values. (Hint: You can use the `getAttributes()` method to detect all available attributes of the current node).

3 JSON

In the second part of this exercise, we focus on the JSON format. As you already have some experience with the *Mondial* dataset in the first step, you are asked to transform parts of the XML into a JSON. In order to do so make use of the Google Gson Java library³.

3.1 Mondial — XML to JSON

TASK 6: Create a JSON file (*.json) which contains all countries which are located in Europe with the attributes of the `country` node from the original `mondial-*.xml`. (Hint: Have a look at the last exercise of the former section. Gson offers a method to simply translate a `HashMap` into a JSON string, which then can be written to a file.)

3.2 Mondial — Reading JSON

In the second step we want to create Java objects from the JSON file we just created, but we are not interested in all attributes.

TASK 7: Write a small program, which reads the JSON file (which was the output of the former task) and transforms each line into a Java object (named `Country.java`). The country should have four values: the `id` (String), the `name` (String), the `car_code` (String), and the `population` (Long). Do you have to pay attention to type conversion? What is the total number of inhabitants of those countries? (Hint: Have a look at the example code of the lecture.)

³You can find the library at the Google code page: <https://sites.google.com/site/gson/>. A user guide can be found on this page: <https://github.com/google/gson/blob/master/UserGuide.md>

4 RDF

In the last part of this exercise session, we will focus on RDF and SPARQL. On the course web page you can find the European countries with their name, population and spoken languages stored as RDF file. The file was generated from the original *mondial* XML file.⁴ In the following you will be asked to formulate SPARQL queries to answer questions about the dataset using the *Jena* Java Framework⁵. In addition to the lecture the *W3* site of *SPARQL Query Language* can help you to answer the questions.⁶

4.1 Mondial — Query with SPARQL I

TASK 8: Write a small program, which reads the RDF file (from the course web page) and formulate a SPARQL query which returns the name and id of all countries within the dataset ordered by the name. What is the last country on this list? In order to explore the property names and namespaces have a look at the RDF file or at the code which was used to generate the file. (Hint: Have a look at the example code of the lecture.)

4.2 Mondial — Query with SPARQL II

As we now have set up the code to query against our dataset, we are interested in the largest countries. But as we already know that Russia and Germany are pretty large, we want to generate a list of the second top 5 largest countries by population.

TASK 9: What is the SPARQL query which returns the second five (6th to 10th) most populated countries in Europe? And which countries are these?

4.3 Mondial — Query with SPARQL III

In the last exercise, you are asked to write a SPARQL query which selects all countries whose inhabitants speak a defined language.

TASK 10: How does the SPARQL query look like, which returns a list of all German-speaking countries with their name and id?

⁴The code which was used to generate the file can also be found in the Java project of this (see `de.dwslab.lecture.wdi.rdf.Converter.java`).

⁵The documentation of the framework can be found at their website: <https://jena.apache.org/>

⁶<https://www.w3.org/TR/rdf-sparql-query/>