

Homework #2: SPARQL and Linked Data

The following exercises are due at the beginning of class on Tuesday, Feb. 19. There are two sections: written exercises and electronic exercises. Pay attention to the extra requirements for CSE 428 students. CSE 398 students can receive extra credit for doing the work expected of CSE 428 students. This will count for 10% of your overall grade.

Written Exercises:

The exercises in this section should be completed and turned in on paper.

1. *[20 pts.]* For this exercise, assume that you are operating on a model using the **swpub.rdf** schema from the last homework. Write the following queries in SPARQL syntax:
 - a) Retrieve the titles of all papers that appear in conference proceedings.
 - b) For all papers written between 2001 and 2005, inclusive, retrieve the title, human-friendly name of the type (i.e., using its `rdfs:label`) and year of the paper.
 - c) Find all papers that appear in a venue with “ISWC” in the name (using partial string matching), and return a list of title, venue, year and topic label (i.e., not the URI of the topic), sorted by year in descending order, followed by topic in ascending order (for those papers in the same year).
 - d) **(Required for CSE 428 students, extra credit for CSE 398 students)** For each author, retrieve their name and a count of the number of papers they have authored. Sort the result by author name. You should assume that two authors are the same if and only if their full names match exactly. For sorting purposes, just use the name as it appears, do not worry about first name or last name distinctions. Hint: You need to use proposed SPARQL 1.1 extensions to do this.
2. *[10 pts.]* Write a SPARQL construct query that generates all triples inferred by RDFS entailment rule `rdfs9` (see the RDF Semantics recommendation [<http://www.w3.org/TR/rdf-mt/>], Section 7). The query should work correctly regardless of domain schema used (i.e., it should be domain independent). You may assume that all entailed all `rdfs:subClassOf` states are already included in the graph.
3. *[10 pts.]* Consider a triple store that contains social networking information in the form of FOAF profiles (<http://www.foaf-project.org/>) collected from various web sites. Assume each pay-level domain uses different URIs for people, but that some sites have users in common. One way to generate links for FOAF is when two resources share the same value for the `foaf:mbox` property, which provides an e-mail address for an agent. Write a SPARQL construct query to generate `owl:sameAs` statements based on this property. Make sure that your query only generates triples where the subject and object are different (i.e., non-trivial `owl:sameAs` statements).
4. *[10 pts.]* **(Required for CSE 428 students, extra credit for CSE 398 students)** Based on your personal experience with how the Internet and Web work, what are the pros and cons of generating links as in #3 above? In your answer, consider both precision (what percentage of the generated links are correct) and recall (what percentage of the needed links are generated by the method).

Electronic Exercise:

The following exercise requires both submission of files and hardcopies of these files. The specified files should be included as attachments to a single e-mail sent to heflin@cse.lehigh.edu with subject line “CSE 398/428 – Homework #2 Submission”.

5. [50 pts.] Using Jena, write a class **MakeReadListHTML** that can read in all the files in a specified directory and create a Web page that lists all of the publications organized by topic. You must read all of the files into a single model and only use SPARQL to retrieve information from the model. From the command-line, your program should run as:

```
java MakeReadListHTML input-directory output-filename
```

The program should read in each file in *input-directory* that has a .rdf suffix into a single Jena model. All required information about the papers must be retrieved using one or more SPARQL queries. The program will create an HTML file called *output-filename* that has a list of papers grouped by topic and sorted by year in descending order. Topics should be output as second-level headings (i.e., with <h2> tags). The format of the output depends on which version of the course you are enrolled in:

CSE 398 Students:

Each paper will be a separate paragraph with title, author list, and year, for example:

Reducing OWL Entailment to Description Logic Satisfiability. Ian Horrocks and Peter Patel-Schneider (2003).

CSE 428 Students:

The paper should have a complete bibliographic description, using all available information in the file. You need to be able to output descriptions of papers of any of the types defined in **swpub.rdf**. Make sure to output papers that are missing optional properties, such as location, pages, publishedBy, publisherLoc, and publishedMonth.

For both sets of students, if the paper has an electronic version, then the title should be hyperlinked to it (using an tag). Make sure that the topic labels appear (instead of the topic URIs) and that the authors are listed in correct order (note, you may assume that each paper will have less than 10 authors, which will make this easier).

As with the last homework, you must use Jena 2.7.4 to parse and query the input files. The Jena distribution includes a number of JAR files, and many (but not all) of these will need to be in your classpath for your program to compile and run. You will need to import classes from the **com.hp.hpl.jena.rdf.model** and **com.hp.hpl.jena.rdf.query** packages. Use the **SwPub.java** file from the last homework to make use of constants for the various classes, properties and other resources defined in **swpub.rdf**.

Your program will be graded on functionality and style. Style includes modularity (avoid repeated code when possible, keep methods under ~60 lines, use multiple classes when appropriate), commenting (all of your files should be reasonably commented, including an initial comment that identifies you as the author and descriptive comments for each class and method), proper indentation, clear names, and use of standard naming conventions. The program should be robust.

Print out your .java file(s) and turn in the hardcopy with the rest of your written answers. Create a zip file *your-user-id-hw2.zip* that contains both your source code (.java) and compiled (.class) files (but do not include any of the Jena files in it). Attach the zip file to the e-mail mentioned above.