

## Information Retrieval

### Practical session n°5: Structured IR at INEX

Create a directory named *practice5*. In this directory, create a new file named *practice5\_report.txt*. For each exercise, copy-paste in this file some outputs of your program, showing that you complete the exercise and it works correctly. Add some explanations.

You will build several runs (maximum = 50 runs / submission) for INEX Ad-Hoc Relevant in Context task. Each run is stored in a textual file, which should be named using the following template (see Practical session n°4 for details): *TeamName\_Run-Id\_WeighingFunction\_Granularity\_Stop\_Stem\_Parameters.txt*

Example: VictorAlbertJulesIsaac\_12\_bm25\_elements\_sec\_p\_stop344\_nostem\_k1.2\_b0.75.txt

At the end of your work:

- copy-paste the source code of your program(s) in the directory *practice5*.
- copy all your runs (maximum = 50 runs) in the directory *practice5* (without any sub-directory).
- compress the directory *practice5* in a file named *practice5\_YourTeamName.zip* (or *.tar*, *.gz*, *.rar*, etc.) (e.g.: *practice5\_VictorAlbertJulesIsaac.zip*).
- upload this compressed file (one file / team) on the website of the course **before 11.59pm November 21<sup>st</sup>**.

The week after, repeat the whole process: create a directory named *practice5v2*, a file *practice5v2\_report.txt*, copy your source code, copy all your runs (maximum = 50 runs), compress the directory *practice5v2* in a file named *practice5v2\_YourTeamName.zip* (or *.tar*, *.gz*, *.rar*, etc.) and upload this compressed file (one file / team) on the website of the course **before 11.59pm December 5<sup>th</sup>**.

And again, repeat the whole process: *practice5v3*, *practice5v3\_report.txt*, *practice5v3\_YourTeamName.zip* ... and upload this compressed file **before 11.59pm December 12<sup>th</sup>**.

And again, repeat the whole process: *practice5v4*, *practice5v4\_report.txt*, *practice5v4\_YourTeamName.zip* ... and upload this compressed file **before 11.59pm January 3<sup>rd</sup>**.

#### Exercise 1: Indexing XML documents

Download the collection of XML documents *Practice\_05\_data.zip* on the website of the course.

It contains 9,804 XML documents stored in 9,804 XML files (237MB).

Index this collection, without considering XML tags, and using your indexing program and the weighting function of your choice.

Compute again the collection statistics and check if they are the same as those computed during the practical session n°3.

#### Exercise 2: SMART *ltn* run (XML elements)

Retrieve a ranked list of 1,500 XML elements from the collection of the Exercise 1, for each of the 7 queries of the Practical session n°4, using SMART *ltn* weighting function, using stop-words or not, and using a stemmer or not. Fix the granularity (i.e. list of XML tags considering as potential indexing unit) as you want. With these results, build one or several runs for the RIC task of INEX. Check carefully the syntax.

#### Exercise 3: SMART *ltc* run (XML elements)

Same question, but using SMART *ltc* weighting function instead of *ltn*.

#### Exercise 4: BM25 run (XML elements)

Same question, but using BM25 weighting function instead of *ltc*. Set the  $k_1$  and  $b$  parameters as you want. You can use these usual values:  $k_1 = 1.2$  and  $b = 0.75$ .

**Exercise 5: fields weighting [Wilkinson94]**

Retrieve a ranked list of 1,500 documents, for each of the 7 queries, using the late combination of fields proposed by [Wilkinson94], and using *BM25* weighting function. Fix the granularity (i.e. list of XML tags considering as fields). Set the  $k_f$ ,  $b$ , and  $\alpha_i$  parameters as you want.

With these results, build one or several runs for the RIC task of INEX.

**Exercise 6: fields weighting [Robertson94]**

Same question, using the early combination of fields proposed by [Robertson94].