

CE706 - Information Retrieval

Assignment 1: Indexing for Web Search

Udo Kruschwitz

5th February 2019

Plagiarism

You are reminded that this work is for credit towards the composite mark in CE706, and that the work you submit must therefore be your own. Any material you make use of, whether it be from textbooks, the Web or any other source must be acknowledged as a comment in the program, and the extent of the reference clearly indicated.

The Context of your Task

Here are some snippets extracted from a job ad for a *Software Engineer* at Microsoft which was posted last Thursday.¹ It starts with the question “Do you want to help make Bing the best search engine in the world?” and then continues with

... We are currently building an entire new system for web search indexing backend that will be orders of magnitude larger and faster than anything that currently exists. The new indexing pipeline is built on top of a petabyte scale table and incremental processing infrastructure. The goal of the system is to process billions of documents a day with seconds to minutes E2E latency. ...

Our core set of engineering challenges include

- ...
- *Parsing and classifying billions of web documents and do encoding & language detection, script segmentation and sentence breaking.*
- *Use NLP techniques to extract Named entities entity linking, POS tagging, stemming etc.*

We write software from the ground-up, running across thousands of servers and managing petabytes of data. Problems our group addresses daily range from

- *designing new infrastructure pieces that lets us scale to handle petabytes of data processing*
- *NLP algorithms to extract features out web documents*
- *Debugging tools for customer feedback on why a particular page does not have the right data in the index.*

That looks exactly like the profile I would like you to have when you graduate!

The Task

Your task is to apply your IR skills to build a processing pipeline that turns a Web site into structured knowledge (thus enhancing your chances of getting the job outlined above). Your system should take HTML pages as input,

¹<https://careers.microsoft.com/us/en/job/580012/Software-Engineer>

process them using the kind of techniques that we have been looking at in the module, and output an index of terms identified in the documents.

This assignment comes in stages. Marks are given for each stage. You may choose not to attempt some stages. You might also implement a system that does not strictly follow the stages but will work in the same way. The stages are as follows:

- **Engineering a Complete System (10%)** The system you develop must be able to read Web pages from a specified set of URLs and produce appropriately formatted output. The Web pages should be processed one at a time using the steps outlined below. The final system should have control over all the individual components so that there is a single call and all the steps outlined below will be performed.
- **HTML Parsing (10%)** Before the text can be analyzed it is necessary to get rid of the HTML tags. The result will be plain text. Note that if you simply delete all HTML tags, you will lose information such as meta tag keywords. Use an appropriate tool to perform this task.
- **Pre-processing: Sentence Splitting, Tokenization and Normalization (10%)** The next step should be to transform the input text into a normal form of your choice.
- **Part-of-Speech Tagging (10%)** The input should be tagged with a suitable part-of-speech tagger, so that the result can then be processed in the next steps.
- **Selecting Keywords and Phrases (20%)** One aim of your system is to identify the words *and* phrases in the text that are most useful for indexing purposes. Your system should remove words which are not useful, such as very frequent words or stopwords, and identify phrases suitable as index terms. Apply *tf.idf* as part of your selection and weighting step.
- **Entity and Relation Extraction (20%)** Apply a named-entity-recogniser (NER tagger) to your text to identify entities. Extract at least person names, locations and organisations. For additional marks also identify relations that hold between these entities.

You will have noticed that the percentages above only add up to 80%. This is because one of the important aspects of the project is that your work should be well documented and your code well commented. **20% of your mark will come from this.** In addition to the actual code you should submit:

- A description of your implementation: what the code does, and the software you used
- Output produced by a run of your system when applied to these two Web pages:
<http://csee.essex.ac.uk/staff/udo/index.html>
<https://www.essex.ac.uk/departments/computer-science-and-electronic-engineering>
- Output produced by *each* stage of the processing pipeline for each of the two files, i.e. in the suggested staged architecture outlined above this would be output produced by the HTML parser, followed by the output of the sentence splitter, the tokenizer etc. Each stage should produce a separate file (however, to calculate weights such as *tf.idf* for the final index terms you will have to consult data derived from *all* documents and it is up to you how exactly you include that step in your processing pipeline).
- A *short* discussion of your solution focussing on functionality implemented and possible improvements and extensions.

You may work in pairs. If you do, you each need to submit the same report (please include information about which two reports should be treated as a pair). Both members of a pair will get the same mark unless there is reason to do otherwise.

You can implement your system either on the Linux or the Windows machines. Python, Java, C/C++, Perl and shell scripts are good choices for this project, but you are by no means restricted to those languages. Identify suitable open-source tools that help you building your pipeline.

Submission

The assignment, which counts for **20%** of the overall mark, should be submitted as a single *zip file* via the electronic submission system by **Friday, 22 February 2019, 11:59 (mid-day)**. *The guidelines about late assignments are explained in the students' handbook.*