

# COMP4650/6490 Document Analysis

## Assignment 1 – IR

In this assignment, your task is to index a document collection into an inverted index, and then measure search performance based on predefined queries.

A new document collection containing more than 10,000 government site descriptions, and a set of predefined queries, is provided for this assignment.

Throughout this assignment:

1. You will develop a better understanding of indexing, including the tokenizer, parser, and normaliser components, and how to improve the search performance given a predefined evaluation metric,
2. You will develop a better understanding of search algorithms, and how to obtain better search results, and
3. You will find the best way to combine an indexer and search algorithm to maximise your performance.

Throughout this assignment you will make changes to the provided code to improve the information retrieval system. In addition, you will produce an answers file with your responses to each question. Your answers file must be a .pdf file named u1234567.pdf where u1234567 is your Uni ID. You should submit both your modified code files and answer.pdf inside a single zip file.

Your answers to coding questions will be marked based on the quality of your code (is it efficient, is it readable, is it extendable, is it correct).

Your answers to discussion questions will be marked based on how convincing your explanations are (are they sufficiently detailed, are they well-reasoned, are they backed by appropriate evidence, are they clear).

### Question 1 – Running queries (20%)

You have been provided with a simple implementation of an inverted index search system, which can be found in files *inverted\_index.py*, *preprocessor.py* and *similarity\_measures.py*. Your first task is to complete the *run\_queries.py* script. The script should

1. Create an inverted index from gov/documents then
2. Run all of the queries from gov/topics and
3. Write the returned results for all of the queries to runs/retrieved.txt in **TREC\_EVAL format**.

Step 1. has already been implemented for you. You will need to edit *inverted\_index.py* to complete the *run\_query* function so that it returns the highest scoring documents. You will then need to complete the *run\_queryies.py* script to implement steps 2. and 3.

The Text Retrieval Conference (TREC) is a popular conference for research on designing better information retrieval systems. We will use their evaluation tool TREC\_EVAL to measure the performance of your system. The tool expects the returned results to be in a file with a very specific format, each line in this file must be of the form:

query\_id Q0 document\_id rank score MY\_IR\_SYSTEM

Where query\_id identifies the query that this document was returned for. Q0 is the literal string Q0. document\_id is the name of the document that was returned. Score is the similarity score between query and document. Rank is the order this document was retrieved in (e.g. a rank of 0 means that this document was retrieved first, rank of 1 means second and so on). MY\_IR\_SYSTEM is the literal string MY\_IR\_SYSTEM. Once you have successfully created the file you should be able to run the *evaluate.py* file to see how well your retrieved documents matched the ground truth relevant documents for the given queries.

In your answers file to this question put a list of your system's scores for each of the evaluation measures, e.g

map: 0.1

Rprec: 0.1

recip\_rank: 0.1

P\_5: 0.1

P\_10: 0.1

P\_15: 0.1

## Question 2 – TF-IDF Similarity (20%)

Currently, the inverted index uses TF similarity to calculate scores. Your next task is to implement TF-IDF similarity functionality. You will need to complete the methods of the TFIDF\_Similarity class in the *similarity\_measures.py* file. You should compute  $\text{TF-IDF}(t, d) = \text{TF}(t, d) * \log(\frac{n}{\text{DF}(t)})$ , where n is the number of documents in the inverted index. You are encouraged to study the corresponding class TF\_Similarity to see how data can be extracted from the inverted index.

After you have successfully implemented TF-IDF Similarity you can change the `--sim` argument of *run\_queries.py* to TFIDF.

Run *evaluate.py* with your new TFIDF\_Similarity index. Write down your new evaluation scores.

### Question 3 – Evaluation measures (20%)

Did your system perform better with TF\_IDF or TF similarity? You will need to decide which of the provided evaluation measures is most useful for evaluating this system. In your answer you should state which measure(s) you are basing your decision on, and why your chosen measure(s) are appropriate specifically for this gov corpus. If you need to make any assumptions about how the system will be used, then make sure to state them.

You should continue using the best performing similarity measure for the next questions.

### Question 4 – Improving the pre-processor (20%)

The pre-processor defines how raw text is converted into tokens to be indexed. Currently, it does whitespace tokenization and porter-stemming. Your task is to change the Preprocessor class to improve your system's evaluation. You may want to try changing the tokenizer, the stemmer, or any implement any other technique mentioned in the lectures. You are strongly encouraged to make use of the NLTK library for this task, documentation can be found here:

<https://www.nltk.org/api/nltk.html>.

Briefly describe **four** different settings of Preprocessor that you tried, and what impact they had on your system's performance. Explain *why* you think that each change increased/decreased your model's evaluation measures. Write down your evaluation measures for the best performing settings that you found.

### Question 5 – Further Modification (20%)

Your final task is to implement some further modification to the system in order to improve its performance. Exactly what kind of modification you implement is left up to you, some examples of modifications you could make:

- Implement BM25 similarity for the system (<http://ipl.cs.aueb.gr/stougiannis/bm25.html>).
- Adjust the inverted\_index to use differently weighted fields.
- Implement some domain specific pre-processing.
- Anything else you can think of.

In your answers describe your best Information Retrieval system, write down its evaluation measures, and explain why you think your changes improved the measures.

Academic Misconduct Policy: All submitted written work and code must be your own (except for any provided starter code, of course) – submitting work other than your own will lead to both a failure on the assignment and a referral of the case to the ANU academic misconduct review procedures:  
ANU Academic Misconduct Procedures