

Lab Class IR:I

By November 12, 2024, solutions for the following exercises have to be submitted: 1, 2, 3, 4, and 5.

Exercise 1 : Architecture

In which conceptual “layer” (i.e., indexing, storage, and retrieval) would you locate the following components or processes of a search engine.

- (a) Crawler
- (b) Language model
- (c) Snippet generation
- (d) Stemming
- (e) Stop word removal
- (f) Query analysis
- (g) Posting list

Exercise 2 : Document processing

Process this document step-by-step according to a standard document processing pipeline:

Information retrieval (IR) in computing and information science is the task of identifying and retrieving information system resources that are relevant to an information need. The information need can be specified in the form of a search query. In the case of document retrieval, queries can be based on full-text or other content-based indexing. Information retrieval is the science [1] of searching for information in a document, searching for documents themselves, and also searching for the metadata that describes data, and for databases of texts, images or sounds.¹

- (a) What is the output after tokenization (ignoring uppercase/lowercase)?
- (b) What is the output after stop word removal (also remove numbers)?
- (c) What is the output after lemmatization?

Exercise 3 : Query processing

Which essential steps happen to a query until it can be used for searching an inverted index?
How do the steps relate to document processing.

Exercise 4 : Posting list

Construct the posting lists of the index terms a , b , and e , given the two documents d_1 and d_2 :

$d_1 = [a, a, b, c, b, b, d, e, b, c, a, e, d, c]$

$d_2 = [d, f, c, b, b, d, a, a, b, c]$

Exercise 5 : Document statistics

Why could it be useful to store the number of documents that a term occurs in?

¹https://en.wikipedia.org/wiki/Information_retrieval