

IR Assignment 2

Homework Submission Guidelines

1. **Due date: 11.05.23 at 23:55**
2. The assignment can be done in pairs
3. Answers can be submitted either in English or Hebrew
4. HW submission should be done via moodle in the corresponding area (by **only** one of the students)
5. Late submission penalty (**5% a day**) for submitting after the assignment's due date
6. Questions / clarifications and more in the dedicated discussion sub-forum in Piazza

Dry part (70%)

Vector space model (15%):

The following matrix represents the word frequencies of four documents d1, d2, d3, d4. Columns represent the documents in the above order; rows represent the vocabulary of six indexed terms a,b,c,d,e,f in that order. (Use **ln.**)

	d1	d2	d3	d4
a	0	1	1	1
b	1	2	0	1
c	2	0	0	0
d	0	0	0	0
e	1	0	1	1
f	7	5	7	2

Assume that the fraction of corpus documents in which each term appears is 10%, 10%, 20%, 5%, 50%, 90% for the terms a, b, c, d, e, and f, respectively.

1. Compute the cosine similarity between d1 and d2 where terms are represented by the tf-idf scheme. (Describe the tf-idf scheme you have used and provide details of the computation. Use **raw tf.**) (5%)
2. Rank the documents in response to the query "a b f". Use the vector space model where document terms are represented by tf and query terms by tf-idf. Provide details of your computations. (Use **raw tf.**) (5%)
3. Show that, for normalized vectors, Euclidean distance and Cosine similarity induce the same document ranking for a given query (5%)

Term Weighting (10%):

1. What causes the short-documents bias effect when using cosine similarity? (5%)
2. Given following weighted tf function:

$$W_{t,d} = \alpha + (1 - \alpha) \frac{tf_{t,d}}{tf_{\max}(d)}$$

Where α is a value between 0 and 1; $tf_{\max}(d)$ is the raw tf of the most frequent term in the document d .

State **one** reason why this weighted tf function is useful and **one** issue that might arise from it. (5%)

Relevance feedback and evaluation (15%)

1. User 'A' submitted a query to a search engine and obtained an ordered result list. Then, the user provided feedback to the engine (4 – the document is highly relevant to the information need expressed by the query, 0 – the document is not relevant)

DocID	Relevance
5	4
2	1
1	1
3	3
4	0

The total number of relevant documents in the collection is 10.
Calculate the **AP**, **precision** and **recall** (at rank 5) (5%)

2. Suggest a version of Rocchio's model that utilizes graded relevance judgments. (5%)
3. Suggest a version of Rocchio's model that utilizes the rank of relevant documents in the list. (5%)

Evaluation (22%)

1. Propose a variant of AP that uses gradual relevance judgments (10%)
2. In what cases evaluation using MAP will yield the same results as evaluation using MRR? Mention at least 4 different cases (6%)
3. Name two different examples where:
 - a. The removal of stopwords reduces the recall. (3%)
 - b. The removal of stopwords reduces precision. (3%)

True/False questions (8%) :

Mark each of the following sentences as true or false and give a short (**but full**) explanation for why your answer is correct:

1. df_t is an inverse measure of the informativeness of term t . (1%)
2. Vector space-based retrieval is always more effective than Boolean retrieval. (2%)
3. In the vector space model, the higher the value of the normalization factor for a document is, the lower are the chances of retrieval for that document. (1%)
4. The stemming process increases the number of unique terms in the index (1%)
5. Values of $\beta > 1$ in F-measure emphasize precision. (1%)
6. In Rocchio's model, q_0 might be closer to the centroid of the relevant documents than q_m . (2%)

Wet part – Intro to Indri (30%)

Part A: (Assignment_2/data/WET_PART_A)

1. The documents to be indexed for Part A are located in the file **docs.txt**
2. Create an Indri index using the following parameters:

```
<parameters>
  <memory>1G</memory>
  <corpus>
    <path> docs.txt path</path>
    <class>trectext</class>
  </corpus>
  <index>Your folder and index name</index>
</parameters>
```

If the index is created correctly you will find a manifest file **inside** the index directory which looks as follows:

```
<corpus>
  <document-base>1</document-base>
  <frequent-terms>0</frequent-terms>
  <maximum-document>5</maximum-
document>
  <total-documents>4</total-documents>
  <total-terms>212</total-terms>
  <unique-terms>140</unique-terms>
</corpus>
```

Run retrieval with the following parameter file:

```
<parameters>
  <memory>1G</memory>
  <index>Path to your index</index>
  <count>5</count>
  <trecFormat>true</trecFormat>
  <baseline>tfidf,k1:1.0,b:0.3</baseline>
</parameters>
```

1. Run a query "corporation" over the collection using the above parameter file
 - a. How many documents did you retrieve?
 - b. How many documents did you expect to retrieve? Perform and explain the change that is needed for getting the additional documents. (Examine the text of documents.)
2. Write a query that will return document D2 first; **use up to 2 words**; explain your choice.
3. Write a query that will return document D1 first; **use up to 2 words**; explain your choice.
4. By running the query: " Michael Jackson" you will retrieve document D4.
 - a. Do you think D4 is relevant to the information need expressed by this query? Explain.
 - b. Type a query for which D4 can be marked as relevant document; **use up to 2 words**; explain (refer to the ranking score assigned to D4 in response of the two queries)

Part B:

1. The files for PartB are located in **Assignment_2/data/WET_PART_B/**
2. In the PartB folder you will find the following files and directories:
 - a. "AP_Coll.tgz" compress file contains AP documents ("database")
 - b. "queries.txt" – query file with 150 queries
 - c. "qrels_AP" file – the AP relevance judgments
 - d. "StopWords.xml" – the INQUERY 418 stopwords list
 - e. "IndriBuildIndex.xml" – build index configuration file
3. Build 2 indexes using the given "database" directory and parameter file "IndriBuildIndex.xml".
 - a. Index 1: **With** stopwords removal and **with** stemming (use "Krovetz" stemmer)
 - b. Index 2: **With** stopwords removal and **without** stemming.**(Note: Create first 2 index directories, one for each index version)**
4. Run retrieval over the 2 created indexes with the following parameter file (using tf.idf weights):

```
<parameters>
  <memory>1G</memory>
  <index>Your index Path</index>
  <count>1000</count>
  < trecFormat>true</ trecFormat>
  <baseline>tfidf,k1:1.0,b:0.3</baseline>
</parameters>
```

5. Use the trec_eval application or any other evaluation toolkit (details can be found in the lecture of week 2) to evaluate the effectiveness of the 2 retrieved lists and complete the following table. Which retrieval result obtained the highest MAP value? Explain.

Stopword Removal	Krovetz Stemmer	MAP	P@5	P@10
With	With			
With	Without			

Submission Instructions:

1. A **PDF** file containing all answers to the questions (Dry and Wet parts).
2. The name of the file as follows:
HW2_Student_1_EMAIL_Student_2_EMAIL.pdf