

DD2476 Search Engines and Information Retrieval Systems

Assignment 2: Ranked Retrieval¹

The purpose of Assignment 2 is to learn how to implement ranked retrieval. You will learn 1) how to include tf-idf scores in the inverted index; 2) how to handle ranked retrieval from multiword queries; 3) how to use PageRank to score documents; and 4) how to combine tf-idf and PageRank scoring.

The recommended reading for Assignment 2 is that of Lectures 4 and 5.

*Assignment 2 is graded, with the requirements for different grades listed below. In the beginning of the oral review session, the teacher will ask you what grade you aim for, and ask questions related to that grade. The assignment can only be presented once (unless you get an F) – you cannot raise your grade by doing additional tasks after the assignment has been examined and given a grade. **Come prepared to the review session!** The review will take 15 minutes or less, so have all papers in order.*

E: Completed Task 2.1-2.5 with some mistakes that could be corrected at the review session.

D: Completed Task 2.1-2.5 without mistakes.

C: E + Completed Task 2.6

B: C + Completed Task 2.7

A: B + Completed Task 2.8

These grades are valid for review March 5, 2018. See the Canvas pages for grading of delayed assignments.

Assignment 2 is intended to take around 50h to complete.

Computing Framework

For Tasks 2.1-2.2, and 2.7, you will be further developing your code from Assignment 1. For Tasks 2.5 and 2.6, you will be using a source code skeleton found in the course directory **/info/DD2476/ir18/lab/pagerank**. (It can also be downloaded from Canvas). Copy this directory to your home directory.

Task 2.1: Ranked Retrieval

Extend the **search** method in the **Searcher** class to implement ranked retrieval. For a given search query, compute the cosine similarity between the tf-idf vector of the query

¹ With contributions by Carl Eriksson, Dmytro Kalpakchi, Jussi Karlsgren, and Hedvig Kjellström.

and the tf-idf-vectors of all matching documents. Then sort documents according to their similarity score.

You will need to add code to the **search** method, so that when this method is called with the **queryType** parameter set to **Index.RANKED_QUERY**, the system should perform ranked retrieval. You will furthermore need to add code to the **PostingsList**, **PostingsEntry**, and **Searcher** classes, to compute the cosine similarity scores of the matching documents. To sort the matching documents, assign the score of each document to the **score** variable in the corresponding **PostingsEntry** object in the postings list returned from the **search** method. If you do this, you can then use the **sort** method in the built-in **java.util.Collections** class.

When you have finished adding to the program, compile and run it, indexing the data set **davisWiki**. Select the "Ranked retrieval" option in the "Search Options" menu, and try the following two search queries:

zombie	attack
which could result in the list	which could result in the list
Found 36 matching document(s)	Found 228 matching document(s)
0. JasonRifkind.f ... 1. Zombie_Walk.f ... 2. EmilyMaas.f ... 3. AliciaEdelman.f ... 4. Kearney_Hall.f ... 5. Spirit_Halloween.f ... 6. Zombies_Reclaim_the_Streets.f ... 7. StevenWong.f ... 8. Measure_Z.f ... 9. Scream.f ... etc.	0. TheWarrior.f ... 1. Measure_Z.f ... 2. Kearney_Hall.f ... 3. Muilop.f ... 4. bg-33p.f ... 5. Furly707.f ... 6. PamAarkes.f ... 7. s.martin.f ... 8. TrustInMe.f ... 9. stevenscott.f ... etc.

With one-word queries, the numbers above are equal to the length-normalized tf_idf scores of each document with respect to the query term. Our lists above were computed with a tf_idf score for document d and query term t :

- $\text{tf_idf}_{dt} = \text{tf}_{dt} * \text{idf}_t / \text{len}_d$
- $\text{idf}_t = \ln(N / \text{df}_t)$

where

- $\text{tf}_{dt} = [\# \text{ occurrences of } t \text{ in } d]$,
- $N = [\# \text{ documents in the corpus}]$,
- $\text{df}_t = [\# \text{ documents in the corpus which contain } t]$,
- $\text{len}_d = [\# \text{ words in } d]$.

The number of documents should be very similar to those listed above. Possible differences may depend on your regular expressions from assignment 1. Depending on exactly how you compute the similarity scores the **ordering of the documents can**

differ somewhat from those produced by your program – this is fine. You can debug your results by manually computing the tf_{dt} and len_d scores for the top ranked documents d for a term t . (The idf score does not influence the ranking since there is only one term.)

There will not be any examination of Task 2.1, it is merely a preparation for Task 2.2.

Task 2.2: Ranked Multiword Retrieval

Modify your program so that it can search for multiword queries, and present a list of ranked matching documents. All documents that include at least one of the search terms should appear in the list of search results.

When you have finished adding to the program, compile and run it, indexing the data set **davisWiki**. Select the "Ranked retrieval" option in the "Search Options" menu, and try the search queries:

zombie attack	money transfer
Found 249 matching document(s)	Found 1599 matching document(s)
0. JasonRifkind.f ... 1. Zombie_Walk.f ... 2. Kearney_Hall.f ... 3. Measure_Z.f ... 4. Spirit_Halloween.f ... 5. EmilyMaas.f ... 6. AliciaEdelman.f ... 7. TheWarrior.f 8. Scream.f ... 9. Zombies_Reclaim_the_Streets.f ... etc.	0. MattLM.f ... 1. Angelique_Tarazi.f ... 2. JordanJohnson.f ... 3. Transfer_Student_Services.f ... 4. NicoleBush.f ... 5. Anthony_Swofford.f ... 6. Title_Companies.f ... 7. Transfer_Student_Association.f ... 8. Munch_Money.f ... 9. money.f ... etc.

Our lists above were computed with the same length-normalized tf_idf scores for each query term as in Task 2.1, weighed together using cosine similarity.

At the review

To pass Task 2.2, you should be able to start the search engine and perform a search in ranked retrieval mode with a query specified by the teacher, that returns the correct number of documents in an order similar to the model solution used by the teachers. You should also be able to explain all parts of the code that you edited.

Task 2.3: Variants of cosine similarity and tf-idf

This is a pen-and-paper task. Consider:

- The query **money transfer**
- The contents of the file **Munch_Money.f**
- The contents of the file **NinadelRosario.f**

What are the lengths of these documents?

	money transfer	Munch_Money.f	NinadelRosario.f
Euclidean	1,4142135624	14,52583904633395	6,557438524302
Manhattan	2	97	37

Use your search engine (how?) to compute the idf of the words **money** and **transfer**:

$$\text{idf}_t = \ln(N/\text{df}_t) \quad \text{money} : 1.0884964320534185 \quad \text{transfer} : 1.7970117542956678$$

What is the term frequency of the words **money** and **transfer** using these variants:

	money Munch_Money.f	transfer Munch_Money.f	money NinadelRosario.f	transfer NinadelRosario.f
tf	5	0	0	1
$1 + \ln(\text{tf})$	1.6989700043360187	0	0	1.0

Consider the vector space where all document are represented by N-dimensional vectors of tf-idf-values for every word in the vocabulary. Now consider the subspace of only contains the dimensions for the words **money** and **transfer**. In this subspace, what are the coordinates of the three documents above?

	money transfer	Munch_Money.f	NinadelRosario.f
$\text{tf} \times \text{idf}$	(1.088, 1.797)	(5.442, 0)	(0, 1.797)
$(1 + \ln(\text{tf})) \times \text{idf}$	(1.088, 1.797)	(1.849, 0)	(0, 1.797)

What is the cosine similarity between the query and the two documents?

	Munch_Money.f	NinadelRosario.f
Euclidean length, tf	0.5181 0.243	0.8557 0.1078
Manhattan length, tf	0.3771 0.0258	0.6229 0.00515

Euclidean length, $1 + \ln(\text{tf})$	0.5181 0.0827	0.8557 0.1078
Manhattan length, $1 + \ln(\text{tf})$	0.3771 0.00875	0.6229 0.00515

What is the cosine similarity if the query coordinates are considered to be (1,1)?

	Munch_Money.f	NinadelRosario.f
Euclidean length, tf	0.7071 0.24339	0.7071 0.10783
Manhattan length, tf	0.5 0.02577	0.5 0.0135
Euclidean length, $1 + \ln(\text{tf})$	0.7071 0.08270	0.7071 0.10783
Manhattan length, $1 + \ln(\text{tf})$	0.5 0.00875	0.5 0.0135

At the review

To pass Task 2.3, fill in the tables above, and be prepared to explain how you computed the numbers. Furthermore, you should be able to explain how the values in the tables would affect the ranking of results for the query **money transfer**.

Task 2.4: What is a good search result?

This task is a continuation of Task 1.5. The purpose is now to assess whether ranked retrieval gives answers with higher precision and recall than unranked intersection retrieval; this will be done on our set of three representative queries.

We first need to learn **how the quality of ranked retrieval results can be measured** – slightly differently from the unranked retrieval in Task 1.5.

Run the program from Task 2.2, indexing the data set **davisWiki**. Select the "Ranked query" option in the "Search options" menu.

You will continue to extend the text file **FirstnameLastname.txt** from Task 1.5, but now with results from the ranked retrieval. Search the indexed data sets with the same query as in Task 1.5:

graduate program mathematics

Inspect the **50 highest ranked** documents. If you already came across that document for the same query in Task 1.5, use the existing relevance label. Otherwise, assess the

relevance of the document for the query. As in Task 1.5, use the following four-point scale:

- (0) Irrelevant document. The document does not contain any information about the topic.
- (1) Marginally relevant document. The document only points to the topic. It does not contain more or other information than the topic description.
- (2) Fairly relevant document. The document contains more information than the topic description but the presentation is not exhaustive.
- (3) Highly relevant document. The document discusses the themes of the topic exhaustively.

[E. Sormunen. Liberal relevance criteria of TREC—Counting on negligible documents? *ACM SIGIR*, 2002]

Add the results into the file from Task 1.5 using the following space-separated format, one line per assessed document:

QUERY_ID DOC_ID RELEVANCE_SCORE

where **QUERY_ID = 1**, **DOC_ID** = the name of the document, **RELEVANCE_SCORE** = [0, 1, 2, 3]. **Do not remove anything from the file, it should contain the union of Task 1.5 and 2.3 document relevance labels.** Like with Task 1.5, send the text file to **jboye@kth.se**.

It should again be noted that there is no objectively correct relevance label for a certain query-document combination! It is a matter of judgment. For difficult cases, write a short note on why you chose the label you did. At the review, you will present three difficult cases.

Plot a **precision-recall graph** for the returned top-50 list, and compute the **precision at 10, 20, 30, 40, and 50** (relevant documents = documents with relevance > 0).

Assume the total number of relevant documents in the corpus to be 100, and estimate the **recall at 10, 20, 30, 40, and 50**.

Compare the precision at 10, 20, 30, 40, 50 for ranked retrieval to the precision for unranked retrieval. *Which precision is the highest? Are there any trends?*

Do the same comparison of recalls. *Which recall is the highest? Is there any correlation between precision at 10, 20, 30, 40, 50, recall at 10, 20, 30, 40, 50?*

At the review

To pass Task 2.4, you should be able to show the text file with labeled documents, in the correct format. You should have emailed it before presenting. You should show the precision-recall graph for the 50 highest ranked documents, and be able to explain the concepts precision-recall graph and precision at K, and give account for these measures for the returned ranked top-50 document list.

You should also be able to discuss the questions in italics.

Task 2.5: Computing PageRank with Power Iteration

The `pagerank` directory contains the file `PageRank.java`, which is compiled simply by

```
javac PageRank.java
```

The program is executed as follows:

```
java -Xmx1g PageRank linkfile
```

for instance

```
java -Xmx1g PageRank links1000.txt
```

The link file `links1000.txt` is also found in the folder `pagerank`. It contains the link structure of a subset of 1000 articles from a 2012 snapshot of Swedish Wikipedia (only links within these 1000 articles have been kept). Each line has the following structure:

```
1;2,3,4,
```

meaning that webpage number **1** is linking to the articles in **2**, **3** and **4**. We are using numbers instead of the actual webpage names for the sake of brevity; however you can translate from numbers to filenames by using the table in `svwikiArticleTitles.txt`.

Note that the numbers in `links1000.txt` and the indices in the `p` array are not the same! Realizing this might save you hours of debugging. Therefore I'm going to write it again:

Note that the numbers in `links1000.txt` and the indices in the `p` array are not the same!

Your task is to **extend the class `PageRank.java` so that it computes the pagerank of a number of Wikipedia articles given their link structure.** Use the standard power iteration method, as described in the lecture notes and in the textbook (Section 21.2.2), and run your program on `links1000.txt`.

Make sure your program prints the pagerank of the 30 highest ranked pages. Use the array `docName` to translate from internal ID numbers to the numbers in the `links1000.txt` file. Compare with the results in the file `svwiki1000_top_30.txt`.

Look up the titles of some documents with high rank, and some documents with low rank. Does the ranking make sense?

At the review

To pass Task 2.5, you should show that the method returns a very similar top-30 ranking for `links1000.txt` to the one given. The difference in rank for a certain document should not be larger than ± 2 positions, and the difference in pagerank value for the documents should not be larger than ± 0.001 .

You should also be able to explain all parts of the code that you edited, and be able to discuss the questions in italics above.

Task 2.6: PageRank with a sparse matrix representation (C)

The representation of the transition matrix as a two-dimensional array is very wasteful, and limits the size of the links graphs the program can handle. In this task you will therefore explore a method to represent the transition matrix by means of hash tables.

Your task is to **extend the class `PageRankSparse.java` so that it computes the pagerank of the davisWiki pages** given their link structure. Use the standard power iteration method, but this time using the sparse hash table representation. Run your program as follows:

```
java -Xmx1g PageRankSparse linksDavis.txt
```

You should be able to process the graph without using more than 1GB of heap space, and it should not take more than 1 minute or so to compute the pageranks. (If it takes more, you should be able to optimize the main loop).

Make sure your program prints the pagerank of the 30 highest ranked pages. Use the array `docName` to translate from internal ID numbers to the numbers in the `linksDavis.txt` file. Compare with the results in the file `davis_top_30.txt`.

Look up the titles of some documents with high rank, and some documents with low rank. Does the ranking make sense?

At the review

To pass Task 2.6, you should show that the method returns a very similar top-30 ranking for `linksDavis.txt` to the one given. The difference in rank for a certain document should not be larger than ± 2 positions, and the difference in pagerank value for the documents should not be larger than ± 0.001 . You should also be able to explain all parts of the code that you wrote.

Task 2.7: Combining tf-idf and PageRank (B)

In this task, you will integrate your results from Task 2.6 into the search engine we have been developing in Assignment 1 and Tasks 2.1-2.2. When doing a ranked query, make sure that the **score is computed as a function of the tf-idf similarity score and the pagerank** of each article in the result set. Design the combined score function so that you can vary the relative effect of tf-idf and pagerank in the scoring.

Use the pageranks you computed from `linksDavis.txt` in Task 2.6. You should pre-compute the pageranks and read them from file at the start of a search engine session.

You will need to add code to the `search` method, so that when this method is called with the `rankingType` parameter set to `Index.TF_IDF`, the system should perform ranked retrieval based on tf-idf score only, with the `rankingType` parameter set to `Index.PAGERANK`, only pagerank should be regarded, and with the `rankingType` parameter set to `Index.COMBINATION`, your combined score function is used to rank the documents.

When your implementation is ready, compile and run it, indexing the data set **davisWiki**. Select the "Ranked retrieval" option in the "Search Options" menu and the "Combination" option in the "Ranking Score" menu, and try the search queries listed in Task 2.2.

Each query should return the same number of matching documents as in Task 2.2. However, the ranking will vary depending on how you use the document pageranks in the score.

What is the effect of letting the tf-idf score dominate this ranking? What is the effect of letting the pagerank dominate? What would be a good strategy for selecting an "optimal" combination? (Remember the quality measures you studied in Task 2.3.)

At the review

To pass Task 2.7, you should present a function for combining tf-idf and pagerank scores where the influence of the two factors can be varied.

You should be able to start the search engine and perform a search in combination, ranked retrieval mode with a query specified by the teacher, that returns the correct number of documents, and be able to discuss the effect of tf-idf and pagerank on the subsequent ranking.

You should also be able to explain all parts of the code that you edited, and be able to discuss the question in italics above.

Task 2.8: Monte-Carlo PageRank Approximation (A)

The task is now to implement the Monte-Carlo methods 1,2,4 and 5 for approximate pagerank computation mentioned in Lecture 5 and in the paper by Avrachenkov et al. listed as course literature.

Run these four variants on **linksDavis.txt**, using $c = 0.85$ and several different settings of N (the number of initiated walks). Compare the four method variants and settings of N in terms of how fast they converge and how similar the solution is to the exact solution. Implement the following goodness measure:

The sum of squared differences between the exact pageranks and the MC-estimated pageranks for the **30 documents with highest exact pagerank** in **linksDavis.txt**.

Plot this goodness measure for all four methods as a function of N .

What do you see? Why do you get this result? Explain and relate to the properties of the (probabilistic) Monte-Carlo methods in contrast to the (deterministic) power iteration method

Do your findings about the difference between the four method variants and the dependence of N support the claims made in the paper by Avrachenkov et al.?

Finally, use your favorite Monte-Carlo method to approximate the pageranks of the full Swedish Wikipedia link structure (in the file `linksSvwiki.txt`). Iterate until the top 30 documents are stable.

At the review

To pass Task 2.8, you should show a record of your experimentation with the four method variants and their N parameter settings for the `linksDavis.txt` graph.

You should be able to discuss the questions in italics, and be able to discuss the differences between the four variants, compare to the claims in the paper, and explain all parts of the code that you wrote.

Finally, show your list of 30 top documents for the `linksSvwiki.txt` graph. Argue for why they are correct by looking up titles of top documents in the file `svwikiArticleTitles.txt`.