

# Web Information Retrieval - Homework 1

Simone Conia, Marco Favorito, Shani Dana Guetta

April 19 2017

## Contents

<b>1</b>	<b>Homework Overview</b>	<b>1</b>
1.1	MG4J . . . . .	1
1.2	Stemmers . . . . .	1
1.3	Scorers . . . . .	1
1.4	The Cranfield Collection . . . . .	2
1.5	The Time Collection . . . . .	2
1.6	Scripts . . . . .	2
<b>2</b>	<b>Using MG4J</b>	<b>3</b>
2.1	Environment setup . . . . .	3
2.2	Creating a collection . . . . .	5
2.3	Building an inverted index . . . . .	5
2.4	Querying an index and ranking the documents . . . . .	6
<b>3</b>	<b>Scores Aggregation</b>	<b>8</b>
3.1	Fagin's Algorithm . . . . .	8
3.2	Threshold Algorithm . . . . .	10
<b>4</b>	<b>Search Engine Evaluation</b>	<b>13</b>
4.1	Average R-Precision . . . . .	13
4.1.1	Computing R-Precision . . . . .	13
4.1.2	Computing Average R-Precision . . . . .	14
4.1.3	Cranfield Collection Results . . . . .	15
4.1.4	Time Collection Results . . . . .	16
4.2	Average nMDCG . . . . .	17
4.2.1	Computing MDCG . . . . .	17
4.2.2	Computing nMDCG . . . . .	18
4.2.3	Computing Average nMDCG . . . . .	18
4.2.4	Cranfield Collection Average nMDCG Results . . . . .	20
4.2.5	Time Collection nMDCG Results . . . . .	25

<b>5</b>	<b>Conclusions</b>	<b>28</b>
5.1	Best configuration . . . . .	28
5.2	Best stemmer . . . . .	30
5.3	Best scorer function . . . . .	32

# 1 Homework Overview

The objectives of this homework are to configure a search engine on two collections of documents, the Cranfield collection and the Time collection, and to find the best configuration in terms of stemming method and scorer function for the search engine.

## 1.1 MG4J

For this homework, the search engine of choice is MG4J (Managing Gigabytes for Java), a free full-text search engine for large document collections written in Java. It allows to index and query large collection of documents with just few simple steps:

- Create a collection of documents
- Create an inverted index on the created collection
- Query the created index

For more details about MG4J visit:  
<http://mg4j.di.unimi.it>

## 1.2 Stemmers

MG4J includes some stemmers, but it also provides support for custom stemmers (like `homework.EnglishStemmerStopwords`). For this homework, the stemmers of interest are the following ones:

- **`it.unimi.di.big.mg4j.index.NullTermProcessor`:** This is the default stemmer in MG4J: a term processor that accepts all terms and does not do any processing.
- **`it.unimi.di.big.mg4j.index.snowball.EnglishStemmer`:** A stemmer for the English language.
- **`homework.EnglishStemmerStopwords`:** A stemmer for the English language able to filter stop words like "a", "an", "the", etc...

## 1.3 Scorer

The scorers taken into considerations are:

- **`CountScorer`:** A trivial scorer function that computes the score by adding the number of occurrences within the current document of each query term.
- **`TfIdfScorer`:** A scorer that implements the TF/IDF ranking formula.

- **BM25Scorer:** A scorer that implements the BM25 ranking scheme. BM25 is the name of a ranking scheme for text derived from the probabilistic model. The essential feature of the scheme is that of assigning to each term appearing in a given document a weight depending both on the count (the number of occurrences of the term in the document), on the frequency (the number of the documents in which the term appears) and on the document length (in words).

## 1.4 The Cranfield Collection

The search engine is evaluated over the Cranfield Collection. This collection consists of:

- 1400 html documents. These documents contain simple science articles.
- 225 predefined queries.
- 1249 relevance assessments, that is, for each query, a set of relevant documents: the Ground Truth.

For more details about the Cranfield collection visit:  
[http://ir.dcs.gla.ac.uk/resources/test\\_collections/cran](http://ir.dcs.gla.ac.uk/resources/test_collections/cran)

## 1.5 The Time Collection

The search engine is evaluated also over the Time Collection. This collection consists of:

- 423 html documents. These documents contain simple news stories.
- 83 predefined queries.
- 324 relevance assessments, that is, for each query, a set of relevant documents: the Ground Truth.

For more details about the Time collection visit:  
[http://ir.dcs.gla.ac.uk/resources/test\\_collections/time/](http://ir.dcs.gla.ac.uk/resources/test_collections/time/)

## 1.6 Scripts

This homework often mentions some bash or python scripts, but only the most relevant ones are presented and discussed in detail. All the scripts can be found at: <https://github.com/MarcoFavorito/wir-hw1>.

## 2 Using MG4J

### 2.1 Environment setup

Before actually using MG4J, the following configuration script is used to set up the environment:

Listing : main.sh

```
1  if [ $# != 4 ]
2      then
3          echo "No argument supplied"
4          echo "Usage: main.sh <dataset_path> <collection_name> <title_weight>
5              <text_weight>"
6          echo -e "Examples:"
7          echo -e "\tmain.sh Cranfield_DATASET cran 2 1"
8          echo -e "\tmain.sh Time_DATASET time 0 1"
9          exit -1
10     fi
11     echo "Setting environment..."
12
13     export dataset_path=$1
14     export collection_name=$2
15     export title_text_weights="$3 $4"
16
17     export HW_DIR="$(pwd)"
18     export script_folder=$HW_DIR/scripts
19     export collection_path=$dataset_path/$collection_name
20     export ground_truth=${collection_path}/${collection_name}_all_queries.tsv
21     export output_path=out/${collection_name}
22
23
24
25     export stemmers=( it.unimi.di.big.mg4j.index.NullTermProcessor it.unimi.
26                       di.big.mg4j.index.snowball.EnglishStemmer homework.
27                       EnglishStemmerStopwords )
28     export stemmer_names=( default english englishsw )
29     export scorer_functions=( CountScorer TfIdfScorer BM25Scorer )
30     export fields=( text title text_and_title )
31     export OUT_ENGLISHW_BM25S=${output_path}/scores/englishsw/BM25Scorer
32
33     export k_values=( 1 3 5 10 )
34
35     echo "Calling the main script..."
36     source $script_folder/run.sh
```

In order to keep everything as clean as possible, the work has been split into smaller tasks. The following script shows how the work has been organized into separate scripts.

Listing : scripts/run.sh

```
1  #!/bin/bash
2
3  debug_msg="MAIN"
4
5  echo
6  echo $debug_msg - start...
7
8  if [ ! -d $script_folder ]; then
9      echo $debug_msg - cannot find \"scripts\" folder...
10     echo $debug_msg - exiting...
11     echo
12     exit 1
13 fi
14
15 source $script_folder/0-clean.sh
16 source $script_folder/1-set-my-classpath.sh
17 source $script_folder/2-create-collection.sh
18 source $script_folder/3-build-index.sh
19 source $script_folder/4-apply-scorer-functions.sh
20 source $script_folder/5-aggregation.sh
21 source $script_folder/6-evaluation.sh
22 source $script_folder/7-make-plots.sh
23
24 rm -R $script_folder/__pycache__
25 rm $script_folder/*.pyc
26
27 echo "MAIN" - done!
28 echo
```

In particular:

- **0-clean.sh** clears any leftover produced by previous executions of the run.sh script.
- **1-set-my-classpath.sh** sets the paths to the MG4J libraries.
- **2-create-collection.sh** creates the document collection to work on.
- **3-build-index.sh** build the inverted indices using various combination of stemmers and scorer functions.
- **4-apply-scorer-functions.sh** queries the indices using a predefined set of queries, and returns the most relevant results.
- **5-aggregation.sh** aggregates the scores of the of some given fields for a given stemmer/scorer configuration.
- **6-evaluation.sh** evaluates the search engines using Average-R-Precision and nMDCG.

- **7-make-plots.sh** draws some plots to show the behavior of the search engine.

## 2.2 Creating a collection

The following script creates a collection from a set of HTML documents. For example, for the Cranfield Collection, it will create a collection named *cran.collection* using all the HTML documents found in the *Cranfield\_DATASET/cran* folder.

Listing : scripts/2-create-collection.sh

```

1  #!/bin/sh
2
3  debug_msg="CREATE_COLLECTION"
4
5  echo
6  echo $debug_msg - start...
7
8  if [ ! -d ${output_path} ]; then
9      echo $debug_msg - creating \"${output_path}\" directory...
10     mkdir ${output_path} --parents
11 fi
12
13 find $collection_path -iname \*.html | \
14     java it.unimi.di.big.mg4j.document.FileSetDocumentCollection \
15         -f HtmlDocumentFactory -p encoding=UTF-8 $output_path/
16         $collection_name.collection
17 echo $debug_msg - done!
18 echo

```

## 2.3 Building an inverted index

The analysis concerns three stemmers: the default one, the English stemmer, and the custom English stemmer that filters stop words. This means that three different indices have to be built. The following script creates three different inverted indices using the three given stemmers, and then it stores them in three separate directories.

Listing : scripts/3-build-index.sh

```

1  #!/bin/sh
2
3  debug_msg="BUILD_INDEX"
4
5  build_index () {
6      stemmer=$1
7      stemmer_name=$2

```

```

8
9     echo $debug_msg - creating new \"${output_path}/indices/${stemmer_name}
    \" folder...
10    mkdir ${output_path}/indices/${stemmer_name} --parents
11
12    java it.unimi.di.big.mg4j.tool.IndexBuilder -t $stemmer -S ${
    output_path}/${collection_name}.collection ${output_path}/indices/
    ${stemmer_name}/${collection_name}
13
14    # create a link for the .collection file into the index folder.
15    ln ${output_path}/${collection_name}.collection ${output_path}/indices/
    ${stemmer_name}/${collection_name}.collection
16  }
17
18  echo
19  echo $debug_msg - start...
20
21  if [ ! -d out ]; then
22      echo $debug_msg - cannot find \"out\" directory...
23      echo $debug_msg - exiting...
24      exit 1
25  fi
26
27  if [ -d ${output_path}/indices ]; then
28      rm -r ${output_path}/indices
29  fi
30  mkdir ${output_path}/indices --parents
31
32  for i in {0..2}; do
33      build_index ${stemmers[i]} ${stemmer_names[i]}
34  done
35
36  echo $debug_msg - done!
37  echo

```

## 2.4 Querying an index and ranking the documents

At this point, a set of predefined queries is given to the search engine. For each query, the search engine will return the most relevant documents and their corresponding scores. These results will be used for score aggregation and search engine evaluation.

---

Listing : scripts/4-apply-scorer-functions.sh

---

```

1  #!/bin/sh
2
3  debug_msg="APPLY_SCORER_FUNCTIONS"
4

```



```

5  echo
6  echo $debug_msg - start
7
8  if [ ! -d ${output_path}/indices ]; then
9      echo $debug_msg - cannot find \"${output_path}/indices\" directory...
10     echo $debug_msg - exiting...
11     exit 1
12 fi
13
14 if [ -d ${output_path}/scores ]; then
15     rm -r ${output_path}/scores
16 fi
17
18 for stemmer in ${stemmer_names[@]}; do
19     for scorer_function in ${scorer_functions[@]}; do
20         mkdir ${output_path}/scores/${stemmer}/${scorer_function} --parents
21
22         for field in ${fields[@]}; do
23             echo "${output_path}/indices/${stemmer}/${collection_name}"
24             java homework.RunAllQueries_HW \
25                 ${output_path}/indices/${stemmer}/${collection_name} \
26                 ${dataset_path}/${collection_name}_all_queries.tsv \
27                 $scorer_function \
28                 $field \
29                 ${output_path}/scores/${stemmer}/${scorer_function}/${
30             field}.tsv
31         done
32     done
33 done
34 echo $debug_msg - done!
35 echo

```

### 3 Scores Aggregation

The following Python scripts implement two scores aggregation algorithms: Fagin's and Threshold. Given a set of rankings, these two algorithms return the top-k elements of a new ranking whose scores are computed by a so called aggregation function.

For this homework, it is requested to aggregate and analyze the scores obtained on the *Title* and *Text* fields when the *BM25Scorer* is applied to the stop-word-filtering-English stemmer. The aggregated score of a document here is defined as the sum of its *Title* score and *Text* score, where the *Title* score is twice as important as the *Text* score.

The output of both the scripts are two files with the following format:

Query_ID	Doc_ID	Rank	Score
...	...	...	...

These files will be used to evaluate the performance of the scores aggregation strategy.

#### 3.1 Fagin's Algorithm

Given a set of rankings, the Fagin's algorithm aggregates these rankings, and then it returns a list containing the top-k document IDs and their corresponding scores.

The Python function takes in input:

- **k** is the number of top documents to return.
- **title\_scorings** is a list of tuples (docID, rank, score), sorted by decreasing score.
- **text\_scorings** is a list of tuples (docID, rank, score), sorted by decreasing score.
- **title\_weight** is the relative importance of the *Title* field.
- **text\_weight** is the relative importance of the *Text* field.

The algorithm is the following one:

1. Sequentially access all lists in parallel until there are k document IDs that have been seen in all lists.
2. Insert in the seen dictionary the document IDs with their score considering the *Title* scores are twice as important as *Text* scores.
3. Perform random accesses to obtain the scores of all the entries that have been seen but not in all the lists.

4. Compute the scores for all objects and return the document IDs with the top-k aggregated scores.

Listing : scripts/fagin.py

```
1  from itertools import zip_longest
2
3  def fagin(k, title_scorings, text_scorings, title_weight, text_weight):
4      """
5      Applies the Fagin's Algorithm on the inputs
6      :param k: the parameter that determines the number of docs returned
7      :param title_scorings: list of tuples (docId, score), sorted by
8                          decreasing score
9      :param text_scorings: list of tuples (docId, score), sorted by
10                         decreasing score
11      :return: a list of tuples (docId, score), with length k, sorted by
12              decreasing score,
13      Notice: the lists @title_scorings and @text_scorings can be of
14              different length (either empty!)
15      """
16      rank = {}
17
18      # get data as dict with key: docId and value: score of that docId for
19      # the query
20      title_dict = dict(title_scorings)
21      text_dict = dict(text_scorings)
22
23      # computes the seen docs, according to the Fagin Algorithm
24      seen = compute_seen_doc_ids(k, title_scorings, text_scorings)
25
26      # compute the score of the seen docs, applying the correct weight
27      # retrieve the scores with "random access" through dictionaries
28      for doc_id in seen:
29          score = title_weight * title_dict.get(doc_id, 0) + text_weight *
30                  text_dict.get(doc_id, 0)
31          rank[doc_id] = score
32
33      # sort the list of (docId, score) by decreasing score and take the
34      # first k
35      sorted_rank = sorted(list(rank.items()), key=lambda x: x[1], reverse=
36                           True)[:k]
37
38      return sorted_rank
39
40      # get data as dict with key: docId and value: score of that docId for
41      # the query
42
43  def compute_seen_doc_ids(k, title_scorings, text_scorings):
```

```

36 # the counter that keeps track of the
37 # number of docId seen in both the lists
38 # until a given iteration
39 counter = 0
40 # the set of seen docIds
41 seen = set()
42
43 # even if the two list of scorings are of variable length,
44 # itertools.zip_longest fills the shorter list with "None" values
45 # and make the two list with equal length
46 for title_entry, text_entry in zip_longest(title_scorings,
47                                           text_scorings):
48     # take the ids if the entry exists, otherwise return None
49     title_entry_id = title_entry[0] if type(title_entry) is tuple else
50     None
51     text_entry_id = text_entry[0] if type(text_entry) is tuple else None
52
53     # if the id is already in the "seen" set, then increase the counter.
54     # otherwise add the current docId in the set.
55     if title_entry_id in seen: counter += 1
56     elif title_entry_id: seen.add(title_entry_id)
57
58     if text_entry_id in seen: counter += 1
59     elif text_entry_id: seen.add(text_entry_id)
60
61     # if we have seen for two times at least k elements,
62     # then stop and return the "seen" set.
63     if counter >= k:
64         break
65
66 return seen

```

## 3.2 Threshold Algorithm

As the Fagin's algorithm, the Threshold Algorithm aggregates these rankings, and then it returns a list containing the top-k document IDs and their corresponding scores.

The Python function takes in input:

- **k** is the number of top documents to return.
- **title\_scorings** is a list of tuples (docID, rank, score), sorted by decreasing score.
- **text\_scorings** is a list of tuples (docID, rank, score), sorted by decreasing score.
- **title\_weight** is the relative importance of the *Title* field.

- **text\_weight** is the relative importance of the *Text* field.

The algorithm is the following one:

1. Consider a row in both lists *title\_scorings* and *text\_scorings*
2. Set the threshold value to be the aggregate of the scores seen in the considered rows.
3. Do random accesses through the dictionaries to compute the score of all seen objects.
4. Sort the scores of the objects seen in the previous step by descending order, and take the top-k objects.
5. If all scores in the top-k list are greater than the threshold value, stop and return this list.
6. Otherwise consider the next row and go back to step 1.

---

Listing : scripts/threshold.py

---

```

1  from itertools import zip_longest
2
3  def threshold(k, title_scorings, text_scorings, title_weight, text_weight
4             ):
5      """
6      Applies the Fagin's Algorithm on the inputs
7      :param k: the parameter that determines the number of docs returned
8      :param title_scorings: list of tuples (docId, score), sorted by
9                          decreasing score
10     :param text_scorings: list of tuples (docId, score), sorted by
11                        decreasing score
12     :return: a list of tuples (docId, rank, score), with length k, sorted
13             by decreasing score,
14     """
15
16     # a dict with key: docId seen and value: score
17     seen = {}
18
19     # get data as dict with key: docId and value: score of that docId for
20     # the query
21     title_dict = dict(title_scorings)
22     text_dict = dict(text_scorings)
23
24     # even if the two list of scorings are of variable length,
25     # itertools.zip_longest fills the shorter list with "None" values
26     # and make the two list with equal length
27     for title_entry, text_entry in zip_longest(title_scorings,
28                                                text_scorings):

```

```

23     # take the pair (ids, score) if the entry exists, otherwise return (
        None, 0)
24     title_entry_id, title_entry_score = (title_entry[0], title_entry[1])
        if type(title_entry) is tuple else (None, 0)
25     text_entry_id, text_entry_score = (text_entry[0], text_entry[1]) if
        type(text_entry) is tuple else (None, 0)
26
27     # compute the threshold of the current iteration
28     cur_threshold = title_weight * title_entry_score + text_weight *
        text_entry_score
29
30     # if the id is not None, add it to the "seen" and associate it with
        its aggregated score
31     # retrieve the scores with "random access" through dictionaries
32     if title_entry_id:
33         seen[title_entry_id] = title_weight * title_dict.get(title_entry_id
            , 0) + text_weight * text_dict.get(title_entry_id, 0)
34     if text_entry_id:
35         seen[text_entry_id] = title_weight * title_dict.get(text_entry_id,
            0) + text_weight * text_dict.get(text_entry_id, 0)
36
37     # sort the list of (docId, score) by decreasing score and take the
        first k
38     top_k = sorted(list(seen.items()), key=lambda x: x[1], reverse=True)
        [:k]
39
40     # IF
41     #   the list of best #k element has length at least #k
42     #   AND
43     #   all the scores of these docs are greater than the current
        threshold
44     # THEN
45     #   RETURN that list
46     if len(top_k) >= k and all([score >= cur_threshold for _, score in
        top_k]):
47         return top_k
48
49     raise Exception("There should be at least k elements...")

```

## 4 Search Engine Evaluation

One of the objectives of this homework is to find the best configuration (in terms of stemming method and scorer function) for the search engine, using the available Ground-Truth data. The following Python scripts implement two metrics used to evaluate the performance of each stemmer/scorer configuration and each aggregation algorithm seen in the previous sections.

### 4.1 Average R-Precision

The first metric used to evaluate the search engine is the Average R-Precision.

#### 4.1.1 Computing R-Precision

Computing R-precision requires knowing a set of documents that are relevant to a query. The number of relevant documents,  $R$ , is used as the cutoff for calculation, and this varies from query to query. For example, if there are 15 documents relevant to "red" in a corpus ( $R = 15$ ), R-precision for "red" looks at the top 15 documents returned, counts the number that are relevant  $r$  turns that into a relevancy fraction:  $r/R = r/15$ <sup>1</sup>.

The Python function that implements R-Precision takes in input:

- **retrieved\_docs** is a list of tuples (docID, rank, score) sorted by decreasing score, representing the ranking computed by the search engine for query of interest.
- **relevant\_docs** is a list that represents the set of relevant documents (docID) specified in the ground-truth file for query of interest.

---

Listing : scripts/r-precision.py

---

```
1 def R_precision(retrieved_docs, relevant_docs):
2     """
3
4     :param retrieved_docs: list of tuples (docId, rank, score), sorted by
        score
5     :param relevant_docs: list of tuples (docId) of the ground-truth
6     :return:
7     """
8
9     relevant_docs_number = len(relevant_docs)
10
11     # take only the best |ground truth| docs, according to the R-Precision
12     top_k_retrieved_docs = retrieved_docs[:relevant_docs_number]
13
14     true_positives = 0
```

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Information\\_retrieval#R-Precision](https://en.wikipedia.org/wiki/Information_retrieval#R-Precision)

```

15 false_positives = 0
16
17 # count the true positives and the true negatives
18 for docId in top_k_retrieved_docs:
19     if docId in relevant_docs:
20         true_positives += 1
21     else:
22         false_positives += 1
23
24 return true_positives / (true_positives + false_positives)

```

#### 4.1.2 Computing Average R-Precision

The Average R-Precision for a set of queries is the mean of the R-Precision scores for each query.

$$\text{Average R-Precision} = \frac{\sum_{q=1}^Q \text{R-Precision}(q)}{Q}$$

where Q is the number of queries.

The Python function that implements Average R-Precision takes in input:

- **all\_retrieved\_docs** is a dictionary with key (queryID) and value a list of tuples (docID, rank, score) sorted by decreasing score.
- **all\_relevant\_docs** is a dictionary with key (queryID) and value a list of documentIDs (docID). This dictionary represents the Ground-Truth.

Listing : scripts/r-precision.py

```

1 def averaged_R_Precision(all_retrieved_docs, all_relevant_docs):
2     """
3     Compute the averaged R-Precision
4     :param all_retrieved_docs: dict[key: qId, value: list[(docId, rank,
5     score)] sorted by score ]
6     :param all_relevant_docs: dict[key: qId, value: list[(docId)]
7     :return: averaged R-Precision among all queries
8     """
9     n_queries = len(all_relevant_docs)
10    precisions = []
11
12    # compute the R-Precision for each query and store in a list
13    for qId in all_relevant_docs.keys():
14        cur_retrieved = all_retrieved_docs.get(qId, [(None, None, None)])
15        cur_retrieved, _, _ = zip(*cur_retrieved)
16        cur_relevant = all_relevant_docs[qId]

```



```

17     cur_R_precision = R_precision(cur_retrieved, cur_relevant)
18     precisions.append(cur_R_precision)
19
20     # compute the averaged R-Precision from the R-Precision of each query
21     average = sum(precisions)/n_queries
22
23     return average

```

### 4.1.3 Cranfield Collection Results

These are the results for the Average R-Precision evaluation on each stemmer-scorer configuration:

Stemmer	Scorer	Field	Average R-Precision
default	CountScorer	text	0.02030036038656728
default	CountScorer	title	0.06077953103815172
default	CountScorer	text_and_title	0.02578652837273527
default	TfIdfScorer	text	0.1943322947633293
default	TfIdfScorer	title	0.17221935023659166
default	TfIdfScorer	text_and_title	0.17932787501753014
default	BM25Scorer	text	0.23548681609026445
default	BM25Scorer	title	0.20043697112662626
default	BM25Scorer	text_and_title	0.2549430381326934
english	CountScorer	text	0.023642988729195624
english	CountScorer	title	0.07690215233318681
english	CountScorer	text_and_title	0.02997251014492394
english	TfIdfScorer	text	0.2083948239120653
english	TfIdfScorer	title	0.17517752560856006
english	TfIdfScorer	text_and_title	0.1893509074543557
english	BM25Scorer	text	0.25755338557062707
english	BM25Scorer	title	0.22468447296033506
english	BM25Scorer	text_and_title	0.26239166583994167
englishsw	CountScorer	text	0.12768598587564103
englishsw	CountScorer	title	0.19686634988359122
englishsw	CountScorer	text_and_title	0.16137814715400925
englishsw	TfIdfScorer	text	0.2081881517226345

englishsw	TfIdfScorer	title	0.17298519626105832
englishsw	TfIdfScorer	text_and_title	0.19088393915980117
englishsw	BM25Scorer	text	0.25153428989635895
englishsw	BM25Scorer	title	0.2213059148403976
englishsw	BM25Scorer	text_and_title	0.2664739776808743

These are the results for the Average R-Precision evaluation on the aggregations algorithms:

Algorithm	Average R-Precision
Fagin	0.24875806987875954
Threshold	0.24875806987875954

#### 4.1.4 Time Collection Results

These are the results for the Average R-Precision evaluation on each stemmer-scorer configuration.

Only the *Text* field is taken into consideration because in the Time Collection the *Title* field has no meaning and it is not informative. This is because the *HtmlDocumentFactory Text* and *Title* fields correspond to the <body> and <head> tags of an Html document. For the Time Collection, the <head> tag contains only the title of the document, so it is not semantically relevant.

Stemmer	Scorer	Field	Average R-Precision
default	CountScorer	text	0.05
default	TfIdfScorer	text	0.0
default	BM25Scorer	text	0.0
english	CountScorer	text	0.03584532248600286
english	TfIdfScorer	text	0.4446586073977995
english	BM25Scorer	text	0.5511641815149965
englishsw	CountScorer	text	0.2392485556022055
englishsw	TfIdfScorer	text	0.4386345110122573
englishsw	BM25Scorer	text	0.5385187947732243

These are the results for the Average R-Precision evaluation on the aggregation algorithms:

Algorithm	Average R-Precision
Fagin	0.5385187947732243
Threshold	0.5385187947732243

In the case of the Time Collection, the aggregation is not effective at all since the *Title* field contains no relevant information, as discussed above. And even if it contained relevant information, just consider the following values, which are the scores of the documents for the *Title* field using the "EnglishStemmerStopWords - BM25Scorer" configuration (the configuration used for the aggregation):

Query_ID	Doc_ID	Rank	Score
13	14	1	2.820362169344532
21	32	1	2.820362169344532
76	31	1	2.820362169344532
76	22	2	2.820362169344532
78	12	1	2.820362169344532

Probably, these scores are so few they do not give enough contribution to make the relatives document IDs important enough to be included in the *top k* documents of the aggregation.

### Some observations

The reader might have noticed that the Average R-Precision of the Fagin's algorithm output and Threshold algorithm output is the same; this result is not by chance but is due to the fact that both Fagin's and Threshold algorithm, if they are used with the same *k*, produce equivalent outputs, i.e. given the same input they return the same output<sup>2</sup>. Obviously, from this follows that the R-Precision evaluation applied over the aggregations is the same, since the aggregation outputs are identical.

## 4.2 Average nMDCG

The evaluation metric nMDCG (normalized Modified Discounted Cumulative Gain) is a modified version of the nDCG (normalized Discounted Cumulative Gain).

### 4.2.1 Computing MDCG

The MDCG of a particular query and value of *k* is defined as follow:

$$\text{MDCG}(\text{query}, k) = \text{relevance}(\text{doc}_1, \text{query}) + \sum_{\text{rank}=2}^k \frac{\text{relevance}(\text{doc}_{\text{rank}}, \text{query})}{\log_2(\text{rank})}$$

where the *relevance* function is defined as follows:

$$\text{relevance}(\text{doc}, \text{query}) = \begin{cases} 1, & \text{doc} \in \text{GroundTruth} . \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

---

<sup>2</sup>The main difference is about performances: Fagin is a simpler algorithm, but at the same time it is not cost optimal for some aggregation functions

### 4.2.2 Computing nMDCG

The nMDCG of a particular query and value of k is defined as follow:

$$\text{nMDCG}(\text{query}, k) = \frac{\text{MDCG}(\text{query}, k)}{\text{MaximumMDCG}(\text{query}, k)}$$

where the *MaximumMDCG* function is defined as follows:

$$\text{MaximumMDCG}(\text{query}, k) = 1 + \sum_{\text{rank}=2}^k \frac{1}{\log_2(\text{rank})}$$

The MaximumMDCG can be seen as the MDCG of a perfect ranking algorithm.

### 4.2.3 Computing Average nMDCG

The Average nMDCG is simply defined as the average value of the nMDCG for each query.

The Python function takes in input:

- **k** that it is requested to be these 4 values 1, 3, 5, 10.
- **all\_retrieved\_docs** is a dictionary [key:queryID, value:list[docId, rank, score]] sorted by score.
- **all\_relevant\_docs** is a dictionary [key:queryID, value:list[docId]]

The functions returns the averaged nMDCG on every query.

---

Listing : scripts/mdcg.py

```
1  from math import log2
2
3
4  def maximumMDCG(k, groundTruth):
5
6      # if k is greater equal the number of relevant docs
7      # (i.e. the length of @groundTruth) then
8      # Compute MDCG with the number of relevant docs.
9      max_k = k if k < len(groundTruth) else len(groundTruth)
10
11     res = 1 + sum([1/log2(rank+1) for rank in range(1, max_k)])
12     return res
13
14 def relevance(docId, groundTruth):
15     return 1 if docId in groundTruth else 0
16
17 def MDCG(k, retrieved_docs, groundTruth):
18     res = 0
19     res += relevance(retrieved_docs[0], groundTruth)
```

```

20     for i in range(1, k):
21         res += relevance(retrieved_docs[i], groundTruth)/log2(i+1)
22
23     return res
24
25
26 def nMDCG(k, retrieved_docs, groundTruth):
27     cur_MDCG = MDCG(k, retrieved_docs, groundTruth)
28     max_MDCG = maximumMDCG(k)
29
30     res = cur_MDCG / max_MDCG
31     return res
32
33
34 def averaged_nMDCG(k, all_retrieved_docs, all_relevant_docs):
35     """
36
37     :param all_retrieved_docs: dict[key: qId, value: list[(docId, rank,
38         score)] sorted by score ]
39     :param all_relevant_docs: dict[key: qId, value: list[(docId)]
40     :return: averaged nMDCG on every queries
41     """
42     nMDCG_values = []
43     queryIds = all_relevant_docs.keys()
44     for qid in queryIds:
45         # if there is no doc for a certain query
46         # we cannot compute nMDCG, so set it to 0
47         if qid not in all_retrieved_docs.keys():
48             cur_nMDCG = 0
49         else:
50             cur_retrieved_docs = all_retrieved_docs[qid]
51
52             # take only the docIds
53             cur_retrieved_doc_ids = [entry[0] for entry in
54                 cur_retrieved_docs]
55             cur_relevant_doc_ids = all_relevant_docs[qid]
56
57             # handle the case when k is greater than the number of
58             # retrieved docs
59             # in that case, k cannot be greater than that number
60             cur_k = k
61             max_k = len(cur_retrieved_docs)
62             if cur_k>max_k:
63                 cur_k=max_k
64
65             # compute nMDCG for the current query
66             cur_nMDCG = nMDCG(cur_k, cur_retrieved_doc_ids,
67                 cur_relevant_doc_ids)
68             nMDCG_values.append(cur_nMDCG)

```

```

65
66     averaged_nMDCG = sum(nMDCG_values)/len(nMDCG_values)
67
68     return averaged_nMDCG

```

#### 4.2.4 Cranfield Collection Average nMDCG Results

These are the nMDCG results for the aggregation algorithms with the different values of K. **Note:** the same considerations stated in Section 4.1.4 holds also in this case.

Algorithm	K	nMDCG
Fagin	1	0.2882882882882883
Threshold	1	0.2882882882882883
Fagin	3	0.30910923604136237
Threshold	3	0.30910923604136237
Fagin	5	0.2870066801783458
Threshold	5	0.2870066801783458
Fagin	10	0.2758719271413992
Threshold	10	0.2758719271413992

The plot in Figure 1 shows a comparison between custom scores aggregation and BM25 Scorer on text and title using the English Stop Words Stemmer.

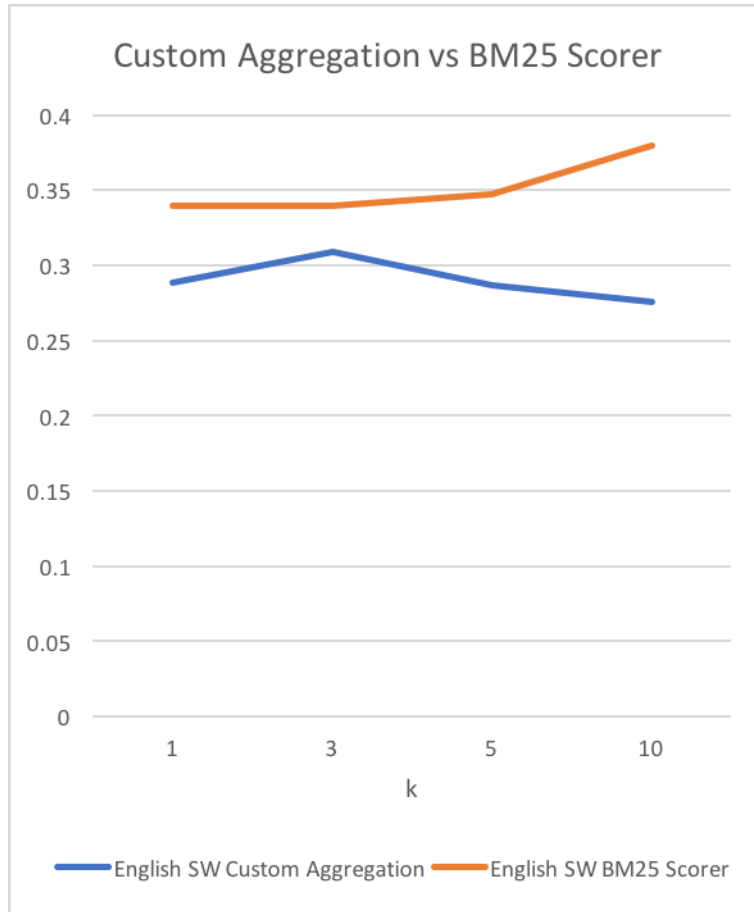


Figure 1: Comparison between custom scores aggregation and BM25 Scorer on text and title using the English Stop Words Stemmer.

Plots in figures 2, 3, and 4 show the Average nMDCG results for the different stemmer-scorer configurations on the Cranfield Collection. On the 'x' axis is represented the value of  $k$  (considering only the following values: 1, 3, 5, 10), while on the 'y' axis the average value of nMDCG over all queries.

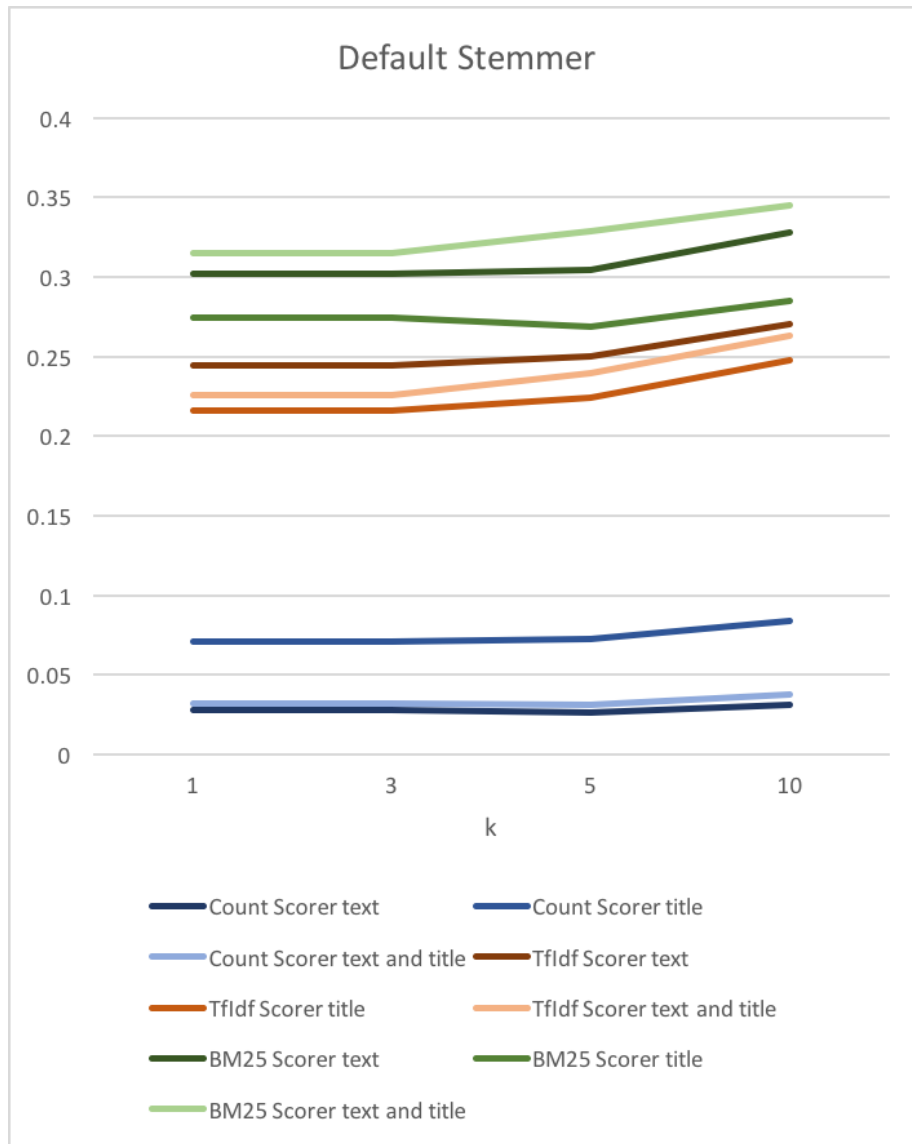


Figure 2: Default Stemmer.



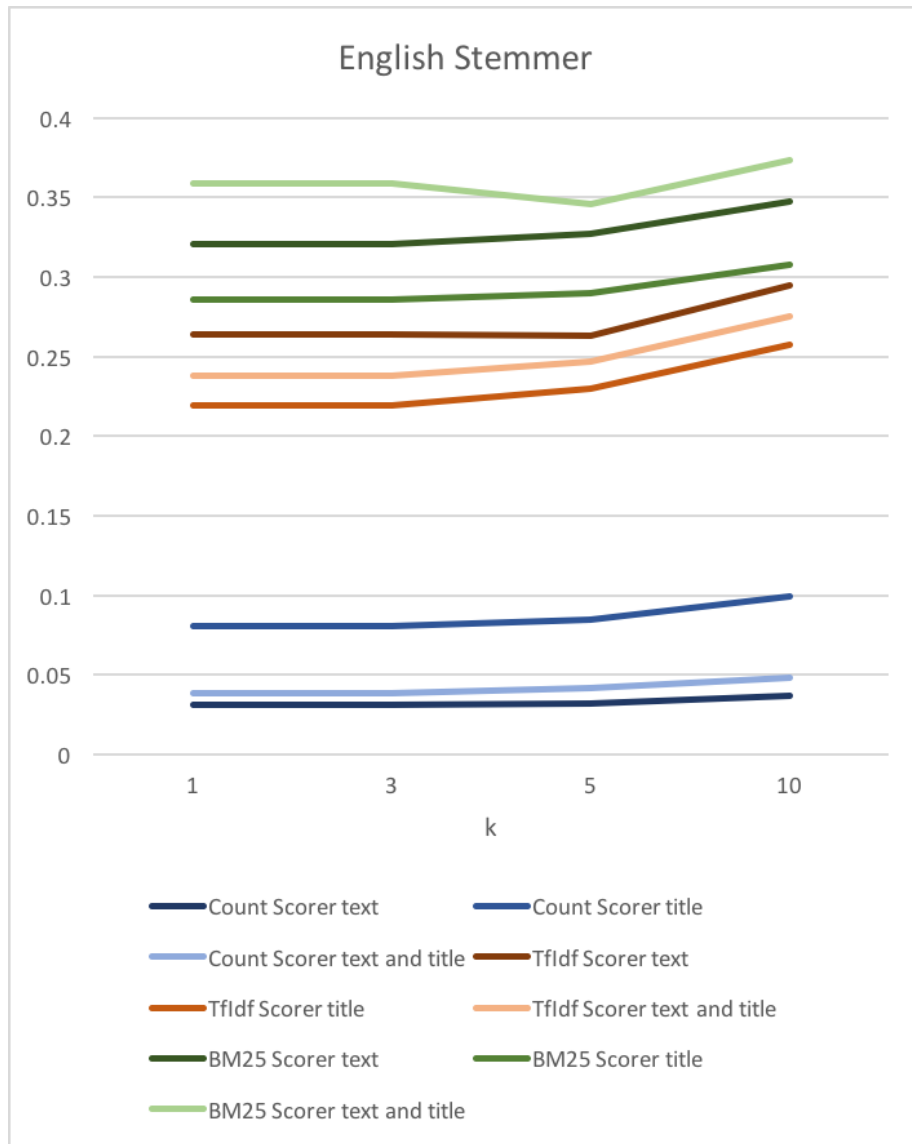


Figure 3: English Stemmer.

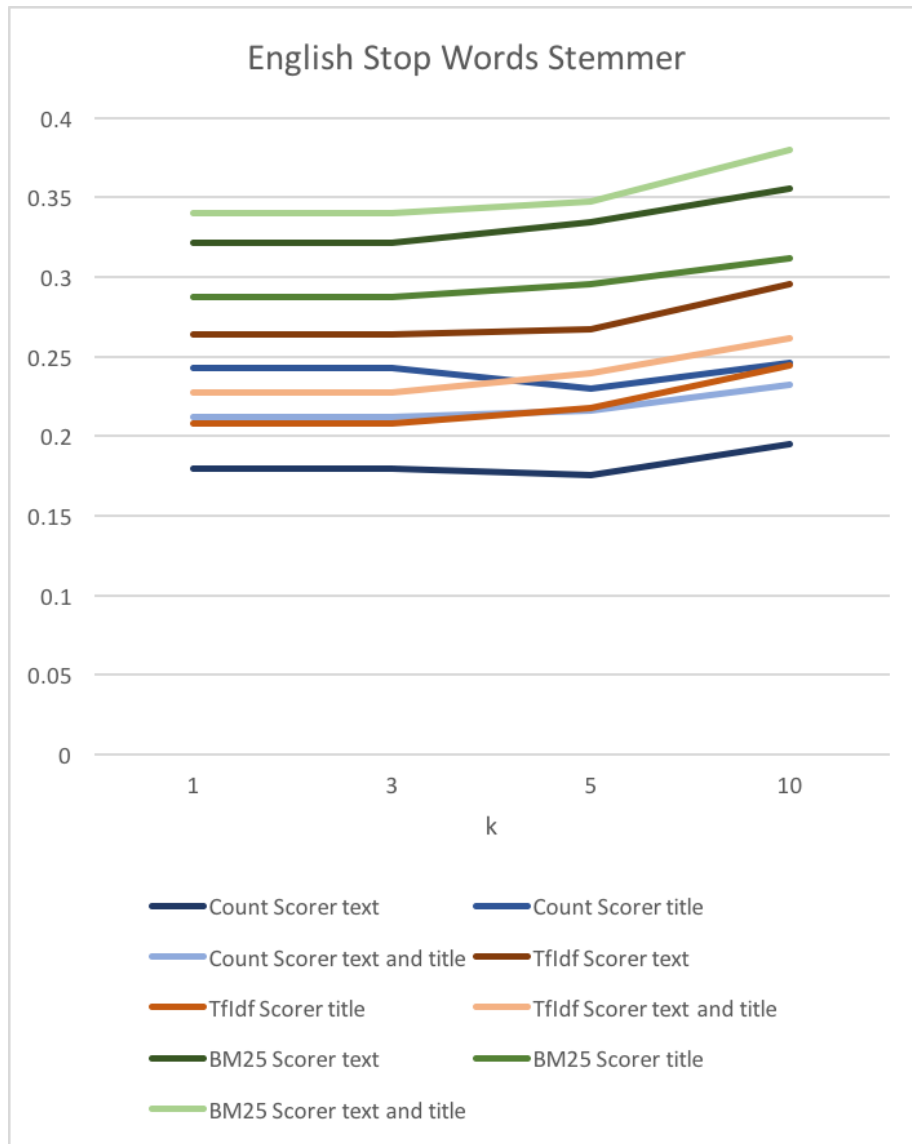


Figure 4: English Stop Words Stemmer.

#### 4.2.5 Time Collection nMDCG Results

These are the nMDCG results for the aggregation algorithms with the different values of K. **Notice:** the same considerations stated in Section 4.1.4 holds also in this case.

Algorithm	K	nMDCG
Fagin	1	0.5903614457831325
Threshold	1	0.5903614457831325
Fagin	3	0.5585508277090627
Threshold	3	0.5585508277090627
Fagin	5	0.5560529519485945
Threshold	5	0.5560529519485945
Fagin	10	0.5474446609447149
Threshold	10	0.5474446609447149

The plot in Figure 5 shows a comparison between custom scores aggregation and BM25 Scorer on text and title using the English Stop Words Stemmer.

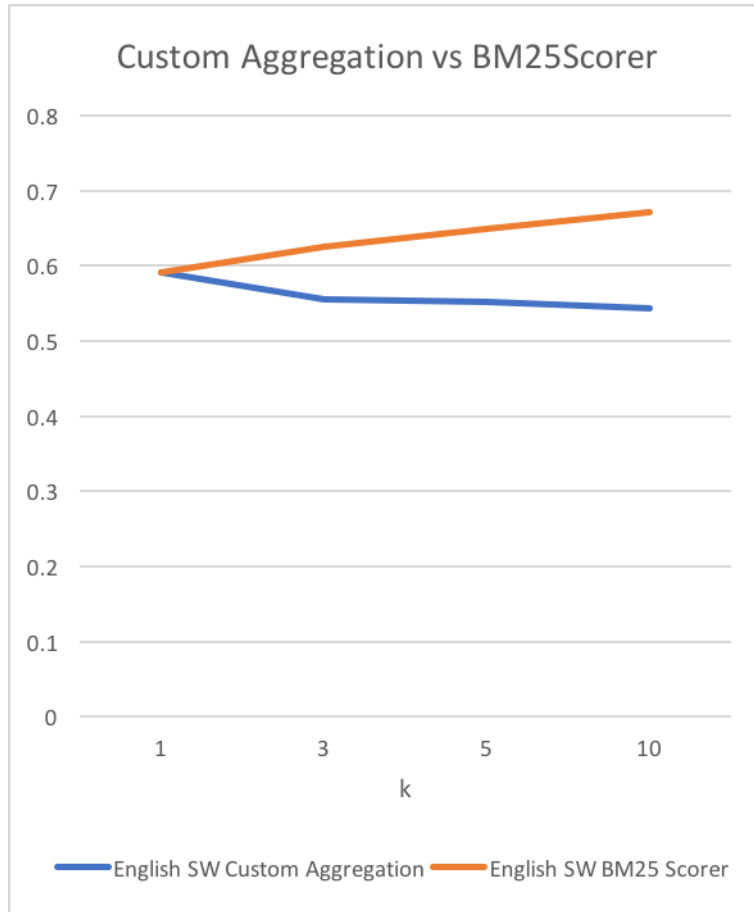


Figure 5: Comparison between custom scores aggregation and BM25 Scorer on text and title using the English Stop Words Stemmer.

Plot in Figure 6 shows the nMDCG results for the different stemmer-scorer configurations on the Time Collection. On the 'x' axis is represented the value of  $k$  (considering only the following values: 1, 3, 5, 10) while on the 'y' axis the average value of nMDCG over all queries.

Only the "Text" field is taken in consideration because in Time Collection the "Title" field has no meaning and it is not informative.

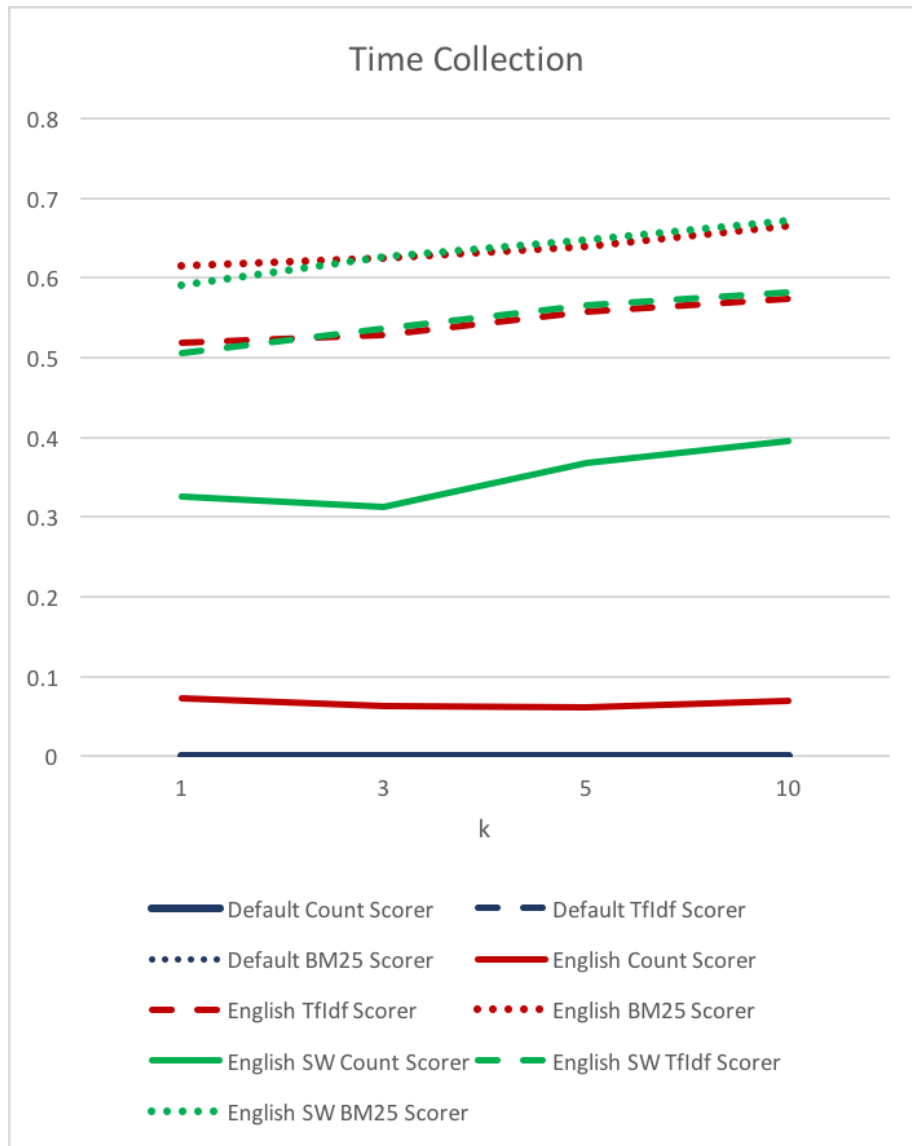


Figure 6: Average nMDCG results on the Text fields for different stemmer/scorer configurations.

## 5 Conclusions

The main purpose of this section is to analyze the previous results in order to answer questions about the best configuration in terms of stemmer and scorer function, the best stemmer, and the best scorer.

**Note:** all the following considerations hold for the scores obtained over both "title" and "text" fields.

### 5.1 Best configuration

The best configuration in terms of stemmer and scorer function is:

- For the Cranfield Collection:
  - According to Averaged R-Precision, "EnglishStemmerStopWords - BM25Scorer" ( $\approx 0.2665$ );
  - According to Averaged nMDCG:
    - \* for  $k = 1$ : all types of stemmer "NullTermProcessor", "EnglishStemmer" and "EnglishStemmerStopWords" with scorer "BM25" ( $\approx 0.3018$ );
    - \* for  $k = 3$ : the best is "EnglishStemmer - BM25Scorer" ( $\approx 0.3586$ ), followed by the pair "EnglishStemmerStopWords - BM25Scorer" ( $\approx 0.3401$ );
    - \* for  $k = 5$ : the best is "EnglishStemmerStopWords - BM25Scorer" ( $\approx 0.3473$ ), followed by the pair "EnglishStemmer - BM25Scorer" ( $\approx 0.3456$ );
    - \* for  $k = 10$ : the best is "EnglishStemmerStopWords - BM25Scorer" ( $\approx 0.3795$ ), followed by the pair "EnglishStemmer - BM25Scorer" ( $\approx 0.3733$ );
- For the Time Collection:
  - According to Averaged R-Precision: "EnglishStemmer - BM25Scorer" ( $\approx 0.5512$ );
  - According to Averaged nMDCG:
    - \* for  $k = 1$ : the best is "EnglishStemmer - BM25Scorer" ( $\approx 0.6144$ ), followed by the pair "EnglishStemmerStopWords - BM25Scorer" ( $\approx 0.6144$ );
    - \* for  $k = 3$ : the best score is reached by "EnglishStemmer - BM25Scorer" ( $\approx 0.6257$ ), followed by "EnglishStemmerStopWords - BM25Scorer" ( $\approx 0.6244$ );
    - \* for  $k = 5$ : the best is "EnglishStemmerStopWords - BM25Scorer" ( $\approx 0.6481$ ), followed by the pair "EnglishStemmer - BM25Scorer" ( $\approx 0.6398$ );

\* for  $k = 10$ : the best is "EnglishStemmerStopWords - BM25Scorer" ( $\approx 0.6717$ ), followed by the pair "EnglishStemmer - BM25Scorer" ( $\approx 0.6648$ );

Considering the above data, the best configuration seems to be either "EnglishStemmer" with BM25 Scorer or "EnglishStemmerStopWords" with BM25 Scorer: they perform almost in the same way as it can be seen in Figure 7.

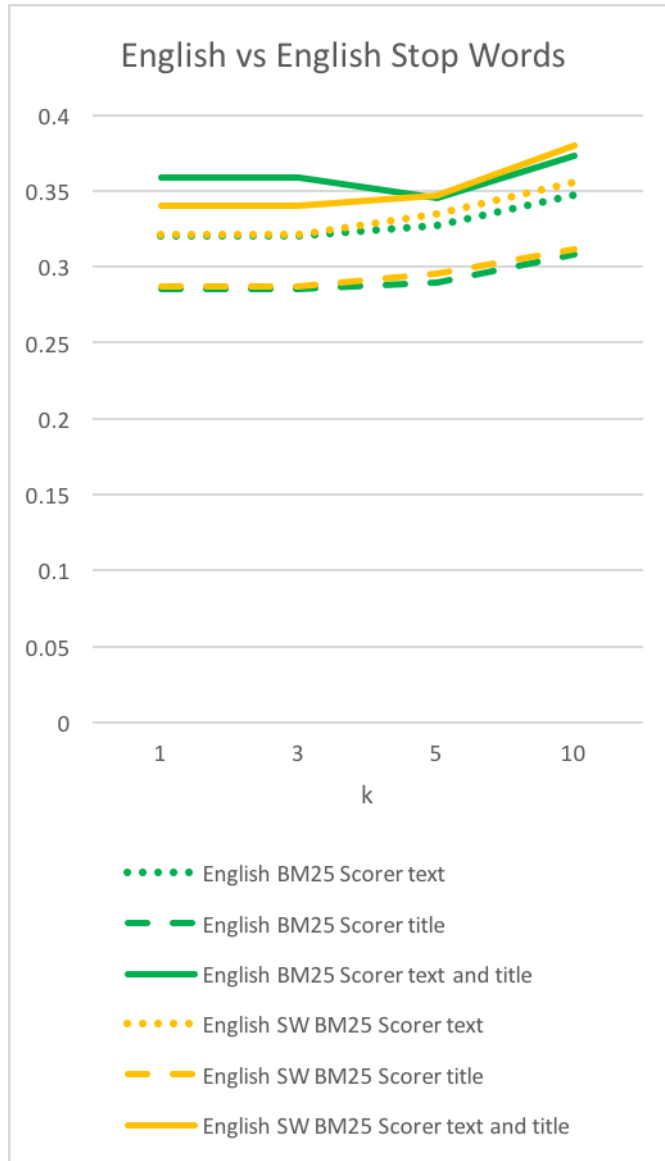


Figure 7: Comparison between English Stemmer and English Stemmer Stop Words both with BM25 Scorer.

## 5.2 Best stemmer

In order to answer to this question, in the following will be reported some useful results computed from the output of the software in `out/${collection_name}/evaluation/benchmark`, both for the Cranfield and the Time collections. In particular, in order to find



the best stemmer with our results, we aggregated the scores by stemmer, averaging over all the scorers.

### **Cranfield Collection**

Evaluating averages for Averaged R-Precision.out

default: 0.153352

english: 0.160572

englishsw: 0.206245

Evaluating averages for Averaged nMDCG with k=1

default: 0.190691

english: 0.201201

englishsw: 0.264264

Evaluating averages for Averaged nMDCG with k=3

default: 0.191278

english: 0.211199

englishsw: 0.260231

Evaluating averages for Averaged nMDCG with k=5

default: 0.200523

english: 0.211233

englishsw: 0.268256

Evaluating averages for Averaged nMDCG with k=10

default: 0.215861

english: 0.231878

englishsw: 0.291918

### **Time Collection**

Evaluating averages for Averaged R-Precision

default: 0.000803

english: 0.343889

englishsw: 0.405467

Evaluating averages for Averaged nMDCG with k=1

default: 0.000000

english: 0.389558

englishsw: 0.457831

Evaluating averages for Averaged nMDCG with k=3

default: 0.000000

english: 0.400708

englishsw: 0.484140

Evaluating averages for Averaged nMDCG with k=5

default: 0.000486

english: 0.417359

englishsw: 0.522976

Evaluating averages for Averaged nMDCG with k=10

default: 0.000486

english: 0.432929

englishsw: 0.549836

### Summary

As we can easily check, in every kind of metric, the best stemmer, on average, is "EnglishStemmerStopWords".

## 5.3 Best scorer function

Doing the same for the scorers (i.e.: aggregating the scores by scorer, averaging over all the stemmers):

### Cranfield Collection

Evaluating averages for Averaged R-Precision

CountScorer: 0.072379

TfIdfScorer: 0.186521

BM25Scorer: 0.261270

Evaluating averages for Averaged nMDCG with k=1

CountScorer: 0.106607

TfIdfScorer: 0.247748

BM25Scorer: 0.301802

Evaluating averages for Averaged nMDCG with k=3

CountScorer: 0.094648

TfIdfScorer: 0.230244

BM25Scorer: 0.337816

Evaluating averages for Averaged nMDCG with k=5

CountScorer: 0.097373

TfIdfScorer: 0.242073

BM25Scorer: 0.340566

Evaluating averages for Averaged nMDCG with k=10

CountScorer: 0.107233

TfIdfScorer: 0.266571

BM25Scorer: 0.365854

## Time Collection

```
Evaluating averages for averaged_r_precision.out  
CountScorer: 0.092501  
TfIdfScorer: 0.294431  
BM25Scorer: 0.363228
```

```
Evaluating averages for averaged_nMDCG_1.out  
CountScorer: 0.136546  
TfIdfScorer: 0.309237  
BM25Scorer: 0.401606
```

```
Evaluating averages for averaged_nMDCG_3.out  
CountScorer: 0.122889  
TfIdfScorer: 0.345231  
BM25Scorer: 0.416728
```

```
Evaluating averages for averaged_nMDCG_5.out  
CountScorer: 0.141518  
TfIdfScorer: 0.370019  
BM25Scorer: 0.429284
```

```
Evaluating averages for averaged_nMDCG_10.out  
CountScorer: 0.157074  
TfIdfScorer: 0.380682  
BM25Scorer: 0.445494
```

## Summary

The best scorer function is BM25, since in every metric, the aggregated score for BM25 is greater than the other scores. It is an expected result, since it is the state-of-the-art TF-IDF-like retrieval functions used in document retrieval<sup>3</sup>.

---

<sup>3</sup>[https://en.wikipedia.org/wiki/Okapi\\_BM25](https://en.wikipedia.org/wiki/Okapi_BM25)