

Homework 1

Ilario Francesco 1469228
Marzio Monticelli 1459333

April 19, 2017

1 Dataset, Stemmers, and Scorer functions

The given Cranfield data-set contains 1400 documents, with a total weight of 1,5 MB, and we have to evaluate 222 queries.

The stemmers we have to use are the following:

- Default Stemmer;
- English Stemmer;
- English Stemmer able to filter Stopwords.

The scorer functions we have to use are the following:

- Count Scorer;
- BM25 Scorer;
- TfIdf Scorer.

2 Software

2.1 Additional Dependencies

- jfreechart: 1.0.0
- jcommon: 1.0.0
- Apache common-CSV: 1.4

2.2 Scripts

2.2.1 Create the collection

This first script, named `1-createCollection.sh`, will create the directory `Collection` and the collection file.

```
1 DIR="./"
2 DIR_COLL=$DIR"Collection/"
3
4 mkdir $DIR_COLL
5
6 find cran -iname \*.html | java it.unimi.di.big.mg4j.document.FileSetDocumentCollection -f
   HtmlDocumentFactory -p encoding=UTF-8 $DIR_COLL"cran.collection"
```

2.2.2 Create the Inverted Indexes

The script `2-createIndexes.sh` creates all the different indexes that we need. It will put them in different directories, every one named after the used Stemmer.

```
1 DIR="./"
2 DIR_COLL=$DIR"Collection/"
3
4 DIR_INDX=$DIR"Indexes/"
5 DIR_INDX_DS=$DIR_INDX"DefaultStemmer/"
6 DIR_INDX_ES=$DIR_INDX"EnglishStemmer/"
7 DIR_INDX_ESS=$DIR_INDX"EnglishStemmerStopwords/"
8
9 mkdir $DIR_INDX
10 mkdir $DIR_INDX_DS
11 mkdir $DIR_INDX_ES
12 mkdir $DIR_INDX_ESS
13
14 cp $DIR_COLL"cran.collection" $DIR_INDX_DS"cran.collection"
15 cp $DIR_COLL"cran.collection" $DIR_INDX_ES"cran.collection"
16 cp $DIR_COLL"cran.collection" $DIR_INDX_ESS"cran.collection"
17
18 java it.unimi.di.big.mg4j.tool.IndexBuilder -S $DIR_COLL"cran.collection" $DIR_INDX_DS"cran"
19
20 java it.unimi.di.big.mg4j.tool.IndexBuilder -t
    it.unimi.di.big.mg4j.index.snowball.EnglishStemmer -S $DIR_COLL"cran.collection"
    $DIR_INDX_ES"cran"
21
22 java it.unimi.di.big.mg4j.tool.IndexBuilder -t homework.EnglishStemmerStopwords -S
    $DIR_COLL"cran.collection" $DIR_INDX_ESS"cran"
```

2.2.3 Obtain the results

In order to obtain the results, we need to execute different scripts. The first is `3-obtainResults.sh`. It evaluates all the given queries -using the `RunAllQueries_HW` software- and saves the results in an ad-hoc output file in an ad-hoc directory. When the script finishes, we find a tree of directories, one for each stemmer, containing three files, one for each “field” between “text”, “title”, and “text_and_title”.

```
1 DIR="./"
2 DIR_RES=$DIR"QueriesResults/"
3 DIR_COLL=$DIR"Collection/"
4
5 DIR_INDX=$DIR"Indexes/"
6 DIR_INDX_DS=$DIR_INDX"DefaultStemmer/"
7 DIR_INDX_ES=$DIR_INDX"EnglishStemmer/"
8 DIR_INDX_ESS=$DIR_INDX"EnglishStemmerStopwords/"
9
10 mkdir $DIR_RES
11
12 ##### Default Stemmer Stopwords #####
13 DIR_DS=$DIR_RES"DefaultStemmer/"
14 DIR_DS_Te=$DIR_DS"DSText/"
15 DIR_DS_Ti=$DIR_DS"DSTitle/"
16 DIR_DS_TT=$DIR_DS"DSTextAndTitle/"
17
18 mkdir $DIR_DS
19 mkdir $DIR_DS_Te
20 mkdir $DIR_DS_Ti
21 mkdir $DIR_DS_TT
22
23 # text
24 java homework.RunAllQueries_HW $DIR_INDX_DS"cran" $DIR"Queries/cran_all_queries.tsv"
    "CountScorer" "text" $DIR_DS_Te"output_text_DefaultStemmer_CountScorer.tsv"
25 java homework.RunAllQueries_HW $DIR_INDX_DS"cran" $DIR"Queries/cran_all_queries.tsv"
    "TfIdfScorer" "text" $DIR_DS_Te"output_text_DefaultStemmer_TfIdfScorer.tsv"
```

```

26 java homework.RunAllQueries_HW $DIR_INDX_DS"cran" $DIR"Queries/cran_all_queries.tsv"
    "BM25Scorer" "text" $DIR_DS_Te"output_text_DefaultStemmer_BM25Scorer.tsv"
27
28 # title
29 java homework.RunAllQueries_HW $DIR_INDX_DS"cran" $DIR"Queries/cran_all_queries.tsv"
    "CountScorer" "title" $DIR_DS_Ti"output_title_DefaultStemmer_CountScorer.tsv"
30 java homework.RunAllQueries_HW $DIR_INDX_DS"cran" $DIR"Queries/cran_all_queries.tsv"
    "TfIdfScorer" "title" $DIR_DS_Ti"output_title_DefaultStemmer_TfIdfScorer.tsv"
31 java homework.RunAllQueries_HW $DIR_INDX_DS"cran" $DIR"Queries/cran_all_queries.tsv"
    "BM25Scorer" "title" $DIR_DS_Ti"output_title_DefaultStemmer_BM25Scorer.tsv"
32
33 # title and text
34 java homework.RunAllQueries_HW $DIR_INDX_DS"cran" $DIR"Queries/cran_all_queries.tsv"
    "CountScorer" "text_and_title"
    $DIR_DS_TT"output_TextAndTitle_DefaultStemmer_CountScorer.tsv"
35 java homework.RunAllQueries_HW $DIR_INDX_DS"cran" $DIR"Queries/cran_all_queries.tsv"
    "TfIdfScorer" "text_and_title"
    $DIR_DS_TT"output_TextAndTitle_DefaultStemmer_TfIdfScorer.tsv"
36 java homework.RunAllQueries_HW $DIR_INDX_DS"cran" $DIR"Queries/cran_all_queries.tsv"
    "BM25Scorer" "text_and_title"
    $DIR_DS_TT"output_TextAndTitle_DefaultStemmer_BM25Scorer.tsv"
37
38 ##### English Stemmer #####
39 DIR_ES=$DIR_RES"EnglishStemmer/"
40 DIR_ES_Te=$DIR_ES"ESText/"
41 DIR_ES_Ti=$DIR_ES"ESTitle/"
42 DIR_ES_TT=$DIR_ES"ESTextAndTitle/"
43
44 mkdir $DIR_ES
45 mkdir $DIR_ES_Te
46 mkdir $DIR_ES_Ti
47 mkdir $DIR_ES_TT
48
49 # text
50 java homework.RunAllQueries_HW $DIR_INDX_ES"cran" $DIR"Queries/cran_all_queries.tsv"
    "CountScorer" "text" $DIR_ES_Te"output_text_EnglishStemmer_CountScorer.tsv"
51 java homework.RunAllQueries_HW $DIR_INDX_ES"cran" $DIR"Queries/cran_all_queries.tsv"
    "TfIdfScorer" "text" $DIR_ES_Te"output_text_EnglishStemmer_TfIdfScorer.tsv"
52 java homework.RunAllQueries_HW $DIR_INDX_ES"cran" $DIR"Queries/cran_all_queries.tsv"
    "BM25Scorer" "text" $DIR_ES_Te"output_text_EnglishStemmer_BM25Scorer.tsv"
53
54 # title
55 java homework.RunAllQueries_HW $DIR_INDX_ES"cran" $DIR"Queries/cran_all_queries.tsv"
    "CountScorer" "title" $DIR_ES_Ti"output_title_EnglishStemmer_CountScorer.tsv"
56 java homework.RunAllQueries_HW $DIR_INDX_ES"cran" $DIR"Queries/cran_all_queries.tsv"
    "TfIdfScorer" "title" $DIR_ES_Ti"output_title_EnglishStemmer_TfIdfScorer.tsv"
57 java homework.RunAllQueries_HW $DIR_INDX_ES"cran" $DIR"Queries/cran_all_queries.tsv"
    "BM25Scorer" "title" $DIR_ES_Ti"output_title_EnglishStemmer_BM25Scorer.tsv"
58
59 # title and text
60 java homework.RunAllQueries_HW $DIR_INDX_ES"cran" $DIR"Queries/cran_all_queries.tsv"
    "CountScorer" "text_and_title"
    $DIR_ES_TT"output_TextAndTitle_EnglishStemmer_CountScorer.tsv"
61 java homework.RunAllQueries_HW $DIR_INDX_ES"cran" $DIR"Queries/cran_all_queries.tsv"
    "TfIdfScorer" "text_and_title"
    $DIR_ES_TT"output_TextAndTitle_EnglishStemmer_TfIdfScorer.tsv"
62 java homework.RunAllQueries_HW $DIR_INDX_ES"cran" $DIR"Queries/cran_all_queries.tsv"
    "BM25Scorer" "text_and_title"
    $DIR_ES_TT"output_TextAndTitle_EnglishStemmer_BM25Scorer.tsv"
63
64 ##### English Stemmer Stopwords #####
65 DIR_ESS=$DIR_RES"EnglishStemmerStopwords/"
66 DIR_ESS_Te=$DIR_ESS"ESSText/"
67 DIR_ESS_Ti=$DIR_ESS"ESSTitle/"
68 DIR_ESS_TT=$DIR_ESS"ESSTextAndTitle/"
69
70 mkdir $DIR_ESS
71 mkdir $DIR_ESS_Te
72 mkdir $DIR_ESS_Ti
73 mkdir $DIR_ESS_TT
74
75 # text
76 java homework.RunAllQueries_HW $DIR_INDX_ESS"cran" $DIR"Queries/cran_all_queries.tsv"
    "CountScorer" "text" $DIR_ESS_Te"output_text_EnglishStemmerStopwords_CountScorer.tsv"

```

```

77 java homework.RunAllQueries_HW $DIR_INDX_ESS"cran" $DIR"Queries/cran_all_queries.tsv"
   "TfIdfScorer" "text" $DIR_ESS_Te"output_text_EnglishStemmerStopwords_TfIdfScorer.tsv"
78 java homework.RunAllQueries_HW $DIR_INDX_ESS"cran" $DIR"Queries/cran_all_queries.tsv"
   "BM25Scorer" "text" $DIR_ESS_Te"output_text_EnglishStemmerStopwords_BM25Scorer.tsv"
79
80 # title
81 java homework.RunAllQueries_HW $DIR_INDX_ESS"cran" $DIR"Queries/cran_all_queries.tsv"
   "CountScorer" "title" $DIR_ESS_Ti"output_title_EnglishStemmerStopwords_CountScorer.tsv"
82 java homework.RunAllQueries_HW $DIR_INDX_ESS"cran" $DIR"Queries/cran_all_queries.tsv"
   "TfIdfScorer" "title" $DIR_ESS_Ti"output_title_EnglishStemmerStopwords_TfIdfScorer.tsv"
83 java homework.RunAllQueries_HW $DIR_INDX_ESS"cran" $DIR"Queries/cran_all_queries.tsv"
   "BM25Scorer" "title" $DIR_ESS_Ti"output_title_EnglishStemmerStopwords_BM25Scorer.tsv"
84
85 # title and text
86 java homework.RunAllQueries_HW $DIR_INDX_ESS"cran" $DIR"Queries/cran_all_queries.tsv"
   "CountScorer" "text_and_title"
   $DIR_ESS_TT"output_TextAndTitle_EnglishStemmerStopwords_CountScorer.tsv"
87 java homework.RunAllQueries_HW $DIR_INDX_ESS"cran" $DIR"Queries/cran_all_queries.tsv"
   "TfIdfScorer" "text_and_title"
   $DIR_ESS_TT"output_TextAndTitle_EnglishStemmerStopwords_TfIdfScorer.tsv"
88 java homework.RunAllQueries_HW $DIR_INDX_ESS"cran" $DIR"Queries/cran_all_queries.tsv"
   "BM25Scorer" "text_and_title"
   $DIR_ESS_TT"output_TextAndTitle_EnglishStemmerStopwords_BM25Scorer.tsv"

```

Subsequently, we can run the script `4-performScoring.sh`. This script calculates the Average R-Precision and the nMDCG on results obtained with the field “text_and_title” for every Stemmer-Scorer function pairs.

```

1  ### scoring ###
2  if [[ $# -eq 1 ]]; then
3      MD=$1
4  elif [[ $# -eq 2 ]]; then
5      MD=$1
6      GTRU=$2
7  else
8      MD="./QueriesResults"
9      GTRU="./Queries/cran_Ground_Truth.tsv"
10 fi
11
12
13 bold=$(tput bold)
14 normal=$(tput sgr0)
15
16 K=(1 3 5 10);
17
18 echo ${bold}"##### AverageRPrecision #####"
19 echo "Default Stemmer - BM25Scorer"${normal}
20 java com.weird.hw1.es.e.AverageRPrecision $MD"/DefaultStemmer/DSTextAndTitle/
   output_TextAndTitle_DefaultStemmer_BM25Scorer.tsv" $GTRU
21
22 echo
23 echo ${bold}"English Stemmer - BM25Scorer"${normal}
24 java com.weird.hw1.es.e.AverageRPrecision $MD"/EnglishStemmer/ESTextAndTitle/
   output_TextAndTitle_EnglishStemmer_BM25Scorer.tsv" $GTRU
25
26 echo
27 echo ${bold}"English Stopwords Stemmer - BM25Scorer"${normal}
28 java com.weird.hw1.es.e.AverageRPrecision $MD"/EnglishStemmerStopwords/ESSTextAndTitle/
   output_TextAndTitle_EnglishStemmerStopwords_BM25Scorer.tsv" $GTRU
29
30 echo
31 echo ${bold}"Default Stemmer - CountScorer"${normal}
32 java com.weird.hw1.es.e.AverageRPrecision $MD"/DefaultStemmer/DSTextAndTitle/
   output_TextAndTitle_DefaultStemmer_CountScorer.tsv" $GTRU
33
34 echo
35 echo ${bold}"English Stemmer - CountScorer"${normal}
36 java com.weird.hw1.es.e.AverageRPrecision $MD"/EnglishStemmer/ESTextAndTitle/
   output_TextAndTitle_EnglishStemmer_CountScorer.tsv" $GTRU
37
38 echo
39 echo ${bold}"English Stopwords Stemmer - CountScorer"${normal}

```

```

40 java com.weird.hw1.esa.AverageRPrecision $MD"/EnglishStemmerStopwords/ESSTextAndTitle/
    output_TextAndTitle_EnglishStemmerStopwords_CountScorer.tsv" $GTRU
41
42 echo
43 echo ${bold}"Default Stemmer - TfIdfScorer"${normal}
44 java com.weird.hw1.esa.AverageRPrecision $MD"/DefaultStemmer/DSTextAndTitle/
    output_TextAndTitle_DefaultStemmer_TfIdfScorer.tsv" $GTRU
45
46 echo
47 echo ${bold}"English Stemmer - TfIdfScorer"${normal}
48 java com.weird.hw1.esa.AverageRPrecision $MD"/EnglishStemmer/ESSTextAndTitle/
    output_TextAndTitle_EnglishStemmer_TfIdfScorer.tsv" $GTRU
49
50 echo
51 echo ${bold}"English Stopwords Stemmer - TfIdfScorer"${normal}
52 java com.weird.hw1.esa.AverageRPrecision $MD"/EnglishStemmerStopwords/ESSTextAndTitle/
    output_TextAndTitle_EnglishStemmerStopwords_TfIdfScorer.tsv" $GTRU
53
54 echo
55 echo ${bold}"##### nMDCG #####"
56 echo "Default Stemmer - BM25Scorer"${normal}
57 for k in ${K[*]}; do
58     java com.weird.hw1.esa.nMDCG -k $k $MD"/DefaultStemmer/DSTextAndTitle/
        output_TextAndTitle_DefaultStemmer_BM25Scorer.tsv" $GTRU
59 done
60
61 echo
62 echo ${bold}"English Stemmer - BM25Scorer"${normal}
63 for k in ${K[*]}; do
64     java com.weird.hw1.esa.nMDCG -k $k $MD"/EnglishStemmer/ESSTextAndTitle/
        output_TextAndTitle_EnglishStemmer_BM25Scorer.tsv" $GTRU
65 done
66
67 echo
68 echo ${bold}"English Stopwords Stemmer - BM25Scorer"${normal}
69 for k in ${K[*]}; do
70     java com.weird.hw1.esa.nMDCG -k $k $MD"/EnglishStemmerStopwords/ESSTextAndTitle/
        output_TextAndTitle_EnglishStemmerStopwords_BM25Scorer.tsv" $GTRU
71 done
72
73 echo
74 echo ${bold}"Default Stemmer - CountScorer"${normal}
75 for k in ${K[*]}; do
76     java com.weird.hw1.esa.nMDCG -k $k $MD"/DefaultStemmer/DSTextAndTitle/
        output_TextAndTitle_DefaultStemmer_CountScorer.tsv" $GTRU
77 done
78
79 echo
80 echo ${bold}"English Stemmer - CountScorer"${normal}
81 for k in ${K[*]}; do
82     java com.weird.hw1.esa.nMDCG -k $k $MD"/EnglishStemmer/ESSTextAndTitle/
        output_TextAndTitle_EnglishStemmer_CountScorer.tsv" $GTRU
83 done
84
85 echo
86 echo ${bold}"English Stopwords Stemmer - CountScorer"${normal}
87 for k in ${K[*]}; do
88     java com.weird.hw1.esa.nMDCG -k $k $MD"/EnglishStemmerStopwords/ESSTextAndTitle/
        output_TextAndTitle_EnglishStemmerStopwords_CountScorer.tsv" $GTRU
89 done
90
91 echo
92 echo ${bold}"Default Stemmer - TfIdfScorer"${normal}
93 for k in ${K[*]}; do
94     java com.weird.hw1.esa.nMDCG -k $k $MD"/DefaultStemmer/DSTextAndTitle/
        output_TextAndTitle_DefaultStemmer_TfIdfScorer.tsv" $GTRU
95 done
96
97 echo
98 echo ${bold}"English Stemmer - TfIdfScorer"${normal}
99 for k in ${K[*]}; do
100     java com.weird.hw1.esa.nMDCG -k $k $MD"/EnglishStemmer/ESSTextAndTitle/
        output_TextAndTitle_EnglishStemmer_TfIdfScorer.tsv" $GTRU
101 done

```

```

102
103 echo
104 echo ${bold}"English Stopwords Stemmer - TfIdfScorer"${normal}
105 for k in ${K[*]}; do
106     java com.weird.hw1.ese.nMDCG -k $k $MD"/EnglishStemmerStopwords/ESSTextAndTitle/
        output_TextAndTitle_EnglishStemmerStopwords_TfIdfScorer.tsv" $GTRU
107 done

```

The third script `-5-ranksAggregation.sh` makes the aggregation of the queries results on separate fields Text, and Title. It uses by default an adaptive value of K (it will be maximum possible for each query), the English Stemmer with Stopwords, and BM25 Scorer function.

```

1 DIR="/"
2 DIR_ESS=$DIR"QueriesResults/EnglishStemmerStopwords/"
3 DIR_ESS_Te=$DIR_ESS"ESSText/"
4 DIR_ESS_Ti=$DIR_ESS"ESSTitle/"
5
6 F_ESS_Te=$DIR_ESS_Te"output_text_EnglishStemmerStopwords_BM25Scorer.tsv"
7 F_ESS_Ti=$DIR_ESS_Ti"output_title_EnglishStemmerStopwords_BM25Scorer.tsv"
8 F_ESS_A=$DIR_ESS"output_aggregated_EnglishStemmerStopwords_BM25Scorer.tsv"
9
10 GTRU="/Queries/cran_Ground_Truth.tsv"
11 R="-r 2"
12 K=""
13 V=""
14
15 while [[ $# -gt 1 ]]
16 do
17     key="$1"
18
19     case $key in
20         -k|-K)
21             K="-k "$2
22             shift
23             ;;
24         -R|-r|--ratio)
25             R="-r "$2
26             shift
27             ;;
28         -t|--text)
29             F_ESS_Te=$2
30             shift
31             ;;
32         -T|--title)
33             F_ESS_Ti=$2
34             shift
35             ;;
36         -o|--output)
37             F_ESS_A=$2
38             shift
39             ;;
40         -v|--verbose)
41             V="--verbose"
42             ;;
43         -V|--Vverbose)
44             V="--Vverbose"
45             ;;
46         *)
47             "usage: "$BASH_SOURCE" [-k|-K K_VALUE -R|--ratio RATIO -t|--text Text_score_FILE
                -T|--title Title_score_FILE -o|--output Output_File]"
48             exit 1
49             ;;
50         esac
51     shift
52 done
53
54 echo
55 # Fagin's Algorithm
56 java com.weird.hw1.rank_aggregation.FaginsAlgorithm $K $R $F_ESS_Te $F_ESS_Ti $F_ESS_AVAG $V
57

```

```

58 echo -ne $bold"Fagin's Algorithm"$normal " "
59 java com.weird.hw1.esa.AverageRPrecision $F_ESS_AVAG $GTRU
60
61 # Fagin's Threshold Algorithm
62 java com.weird.hw1.rank_aggregation.ThresholdAlgorithm2 $K $R $F_ESS_Te $F_ESS_Ti
   $F_ESS_ATHR $V
63
64 echo -ne $bold"Fagin's Threshold Algorithm"$normal " "
65 java com.weird.hw1.esa.AverageRPrecision $F_ESS_ATHR $GTRU

```

At the end, it is possible to plot the nMDCG results using the following script (6-plot.sh):

```

1 # Constants
2 DIROUT="./Plots/"
3 mkdir $DIROUT
4
5 DIR_M="."
6 DIR_QRES=$DIR_M"QueriesResults/"
7 DIR_RES=$DIR_M"Queries/"
8
9 DIR_DS_TT=$DIR_QRES"DefaultStemmer/DSTextAndTitle/"
10 DIR_ES_TT=$DIR_QRES"EnglishStemmer/ESTextAndTitle/"
11 DIR_ES_TT=$DIR_QRES"EnglishStemmerStopwords/ESSTextAndTitle/"
12 DSBM25=$DIR_DS_TT"output_TextAndTitle_DefaultStemmer_BM25Scorer.tsv"
13 DSCOUN=$DIR_DS_TT"output_TextAndTitle_DefaultStemmer_CountScorer.tsv"
14 DSTFIF=$DIR_DS_TT"output_TextAndTitle_DefaultStemmer_TfIdfScorer.tsv"
15
16 ESBM25=$DIR_ES_TT"output_TextAndTitle_EngishStemmer_BM25Scorer.tsv"
17 ESCOUN=$DIR_ES_TT"output_TextAndTitle_EngishStemmer_CountScorer.tsv"
18 ESTFIF=$DIR_ES_TT"output_TextAndTitle_EngishStemmer_TfIdfScorer.tsv"
19
20 ESSBM25=$DIR_ES_TT"output_TextAndTitle_EngishStemmerStopwords_BM25Scorer.tsv"
21 ESSCOUN=$DIR_ES_TT"output_TextAndTitle_EngishStemmerStopwords_CountScorer.tsv"
22 ESSTFIF=$DIR_ES_TT"output_TextAndTitle_EngishStemmerStopwords_TfIdfScorer.tsv"
23
24 GTRU=$DIR_RES"cran_Ground_Truth.tsv"
25
26
27 K="1 3 5 10"
28
29 bold=$(tput bold)
30 normal=$(tput sgr0)
31
32 # Bar Chart
33 echo
34 echo $bold"##### BarChart #####"$normal
35 java com.weird.hw1.plot.BarPlotHwMDCGResults -K $K -f "$DSBM25" "Default Stemmer" "BM25"
   "$DSCOUN" "Default Stemmer" "Count" "$DSTFIF" "Default Stemmer" "TfIdf" "$ESBM25"
   "English Stemmer" "BM25" "$ESCOUN" "English Stemmer" "Count" "$ESTFIF" "English
   Stemmer" "TfIdf" "$ESSBM25" "English Stopwords Stemmer" "BM25" "$ESSCOUN" "English
   Stopwords Stemmer" "Count" "$ESSTFIF" "English Stopwords Stemmer" "TfIdf" -g "$GTRU"
   -o $DIROUT"barChart_All.jpeg"
36
37 # Curves
38 # Default Stemmer
39 echo $bold"##### Curves #####"$normal
40 echo
41 echo $bold"##### Default Stemmer #####"$normal
42 java com.weird.hw1.plot.CurvePlotternMDCG -K $K -e "$DSBM25" -g "$GTRU" -o
   $DIROUT"plot_DefaultStemmer_BM25Scorer.jpeg" -T "Default Stemmer - BM25" -c "blue"
43 echo
44 java com.weird.hw1.plot.CurvePlotternMDCG -K $K -e "$DSCOUN" -g "$GTRU" -o
   $DIROUT"plot_DefaultStemmer_CountScorer.jpeg" -T "Default Stemmer - Count" -c "red"
45 echo
46 java com.weird.hw1.plot.CurvePlotternMDCG -K $K -e "$DSTFIF" -g "$GTRU" -o
   $DIROUT"plot_DefaultStemmer_TfIdfScorer.jpeg" -T "Default Stemmer - TfIdf" -c "green"
47
48 # English Stemmer
49 echo
50 echo $bold"##### English Stemmer #####"$normal
51 java com.weird.hw1.plot.CurvePlotternMDCG -K $K -e "$ESBM25" -g "$GTRU" -o

```

```

$DIROUT"plot_EnglishStemmer_BM25Scorer.jpeg" -T "English Stemmer - BM25" -c "blue"
52 echo
53 java com.weird.hw1.plot.CurvePlotternMDCG -K $K -e "$ESCOUN" -g "$GTRU" -o
$DIROUT"plot_EnglishStemmer_CountScorer.jpeg" -T "English Stemmer - Count" -c "red"
54 echo
55 java com.weird.hw1.plot.CurvePlotternMDCG -K $K -e "$ESTFIF" -g "$GTRU" -o
$DIROUT"plot_EnglishStemmer_TfIdfScorer.jpeg" -T "English Stemmer - TfIdf" -c "green"
56
57 # English Stemmer Stopwords
58 echo
59 echo $bold"##### English Stemmer with Stopwords #####"$normal
60 java com.weird.hw1.plot.CurvePlotternMDCG -K $K -e "$ESSBM25" -g "$GTRU" -o
$DIROUT"plot_EnglishStemmerStopwords_BM25Scorer.jpeg" -T "English Stemmer with
Stopwords - BM25" -c "blue"
61 echo
62 java com.weird.hw1.plot.CurvePlotternMDCG -K $K -e "$ESSCOUN" -g "$GTRU" -o
$DIROUT"plot_EnglishStemmerStopwords_CountScorer.jpeg" -T "English Stemmer with
Stopwords - Count" -c "red"
63 echo
64 java com.weird.hw1.plot.CurvePlotternMDCG -K $K -e "$ESSTFIF" -g "$GTRU" -o
$DIROUT"plot_EnglishStemmerStopwords_TfIdfScorer.jpeg" -T "English Stemmer with
Stopwords - TfIdf" -c "green"

```

2.3 Little Tip

We suggest to use the script 0-executeAll.sh in order to run all the previous ones. It asks to you to tune some parameters or use the default values.

```

1 K_DEF=5
2 re='^[0-9]+$'
3 V=""
4
5 NUM_SCRIPTS=6
6
7
8 bold=$(tput bold)
9 normal=$(tput sgr0)
10
11 # Preparing Help
12 if [[ $# -ge 1 ]]; then
13     for arg in $@; do
14         if [[ $arg == "--help" || $arg == "-h" ]]; then
15             me=$BASH_SOURCE
16             echo "$ source" $me:"
17             echo " -h --help prints this help"
18             echo " -v --verbose sets the verbose mode"
19         fi
20     done
21 fi
22
23 # Handle verbose mode
24 if [[ $# -ne 1 || ( $1 != "-v" && $1 != "--verbose" ) ]]; then
25     V=" > /dev/null"
26     echo ${bold}"[1/"$NUM_SCRIPTS]"${normal}" Performing collection creation"
27 fi
28
29 eval source 1-createCollection.sh $V
30
31 if [[ $V != "" ]]; then
32     echo ${bold}"[2/"$NUM_SCRIPTS]"${normal}" Performing Indexes creation"
33 fi
34 eval source 2-createIndexes.sh $V
35
36 if [[ $V != "" ]]; then
37     echo ${bold}"[3/"$NUM_SCRIPTS]"${normal}" Obtaining Results"
38 fi
39
40 eval source 3-obtainResults.sh $V
41
42 echo

```



```

43 echo ${bold}"[4/"$NUM_SCRIPTS"]"${normal}" Press Enter to perform scoring"
44 read
45 clear
46
47 source 4-performScoring.sh
48
49 # grub K Value
50 echo
51 echo ${bold}"[5/"$NUM_SCRIPTS"]"${normal}" Rank Aggreagation"
52 while (true); do
53     echo -n "which value for K do you want to use? [Press ENTER to use the max possible K
        value for each instance] "
54     read K
55     if [[ $K =~ $re ]]; then
56         K=-k $K
57         break;
58     elif [[ $K == "" ]]; then
59         K=""
60         break;
61     else
62         echo "[ERR] \"$K\" is not a number. Retry" >&2;
63     fi
64 done
65
66 # grub Ratio value
67 while (true); do
68     echo -n "which value for Ratio (importance of Title wrt Text) do you want to use?
        [Default=2] "
69     read R
70     if [[ $R =~ $re ]]; then
71         R=-r $R
72         break;
73     elif [[ $R == "" ]]; then
74         R=""
75         break;
76     else
77         echo "[ERR] \"$R\" is not a number. Retry" >&2;
78     fi
79 done
80
81 eval source 5-ranksAggregation.sh $K $R
82
83 echo
84 echo ${bold}"[6/"$NUM_SCRIPTS"]"${normal}" Plot Results"
85 eval source 6-plot.sh

```

3 Results

3.1 Average R-Precision

	Default	English	English With Stopwords
BM25	0.255	0.262	0.266
Count	0.023	0.030	0.160
TfIdf	0.179	0.189	0.191

Table 1: Average R-Precision for Scorers

	English With Stopwords - BM25
Fagin	0.249
Threshold	0.249

Table 2: Average R-Precision for Ranks Aggregations

3.2 More Rank Aggregations Results

The following results are obtained using the English Stemmer with Stopwords and the BM25 Scorer Function.

	1	3	5	10	20	∞
Fagin	0.067	0.178	0.215	0.243	0.248	0.249
Threshold	0.067	0.178	0.215	0.243	0.248	0.249

Table 3: More Average R-Precision for Ranks Aggregations

4 Plots

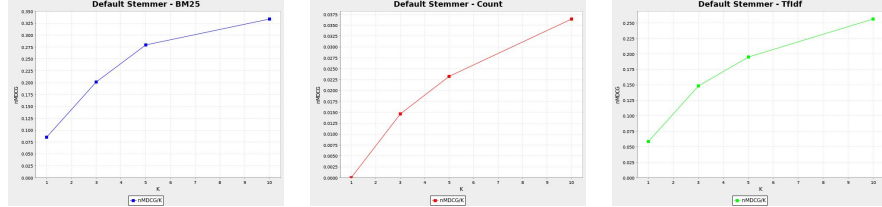


Figure 1: Default Stemmer

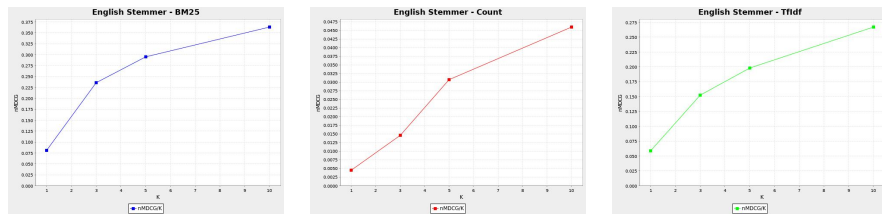


Figure 2: English Stemmer

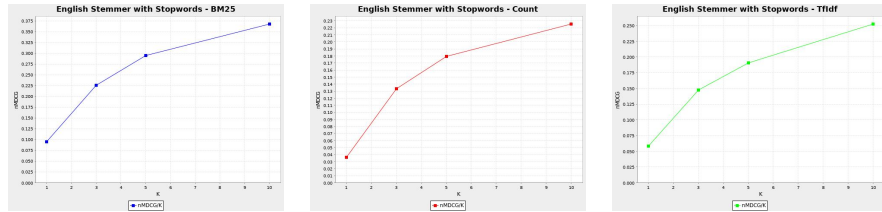
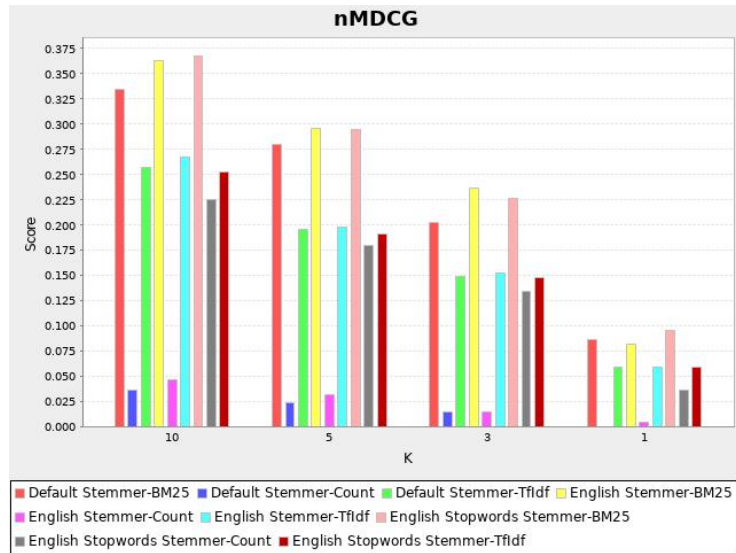


Figure 3: English Stemmer Stopwords



5 Answers

5.1 Best Stemmer-Scorer Function combination

	Default	English	English With Stopwords
BM25	0.255	0.262	0.266
Count	0.022	0.028	0.159
TfIdf	0.179	0.189	0.191

Table 4: (Scorer Functions / Stemmer) Using Average RPrecision

In order to extract this data the K value has been set to 0. Doing this, the software will always use the greatest K possible.

	Default	English	English With Stopwords
BM25	0.406	0.438	0.443
Count	0.058	0.069	0.295
TfIdf	0.318	0.331	0.326

Table 5: (Scorer Functions / Stemmer) Using nMDCG

From tables 4 and 5, we can notice that the best Stemmer-Scorer Function combination is **English Stemmer with Stopwords** and **BM25**.

5.2 Best Stemmer

	Default	English	English With Stopwords
BM25	0.255	0.262	0.266
Count	0.022	0.028	0.159
TfIdf	0.179	0.189	0.191

Table 6: (Scorer Functions / Stemmer) Using Average RPrecision

In order to extract this data the K value has been set to 0. Doing this, the software will always use the greatest K possible.

	Default	English	English With Stopwords
BM25	0.406	0.438	0.443
Count	0.058	0.069	0.295
TfIdf	0.318	0.331	0.326

Table 7: (Scorer Functions / Stemmer) Using nMDCG

From tables 6 and 7, we can notice that the best stemmer is **English Stemmer with Stopwords**.

5.3 Best Scorer Function

	Default	English	English With Stopwords
BM25	0.255	0.262	0.266
Count	0.022	0.028	0.159
TfIdf	0.179	0.189	0.191

Table 8: (Scorer Functions / Stemmer) Using Average RPrecision

In order to extract this data the K value has been set to 0. Doing this, the software will always use the greatest K possible.

	Default	English	English With Stopwords
BM25	0.406	0.438	0.443
Count	0.058	0.069	0.295
TfIdf	0.318	0.331	0.326

Table 9: (Scorer Functions / Stemmer) Using nMDCG

From tables 8 and 9, we can notice that the best Scorer Function is **BM25**.