

# **Assignment 4**

**Fall 2016**

**CS834 Introduction to Information Retrieval**

**Dr. Michael Nelson**

Mathew Chaney

December 4, 2016

## Contents

<b>1</b>	<b>Question 8.3</b>	<b>3</b>
1.1	Question . . . . .	3
1.2	Approach . . . . .	3
1.3	Results . . . . .	4
<b>2</b>	<b>Question 8.4</b>	<b>5</b>
2.1	Question . . . . .	5
2.2	Approach . . . . .	5
<b>3</b>	<b>Appendix</b>	<b>7</b>
3.1	Code listings . . . . .	7
<b>4</b>	<b>References</b>	<b>11</b>

## Listings

1	Output from running the getrel.py script for queries 1 and 10 from the CACM collection.	4
2	The run function . . . . .	5
3	getrel.py . . . . .	7
4	q83.py . . . . .	9
5	q84.py . . . . .	9
6	Script used to generate the recall-precision graphs . . . . .	10

## List of Figures

1	Uninterpolated Recall-Precision Graph for CACM Queries 6 and 8. . . . .	5
2	Average Interpolated Recall-Precision Graph for CACM Queries 6 and 8. . . . .	6

## List of Tables

1	Calculations for CACM query 10 from top 1000 retrieved documents. . . . .	4
---	---	---

# 1 Question 8.3

## 1.1 Question

For one query in the CACM collection (provided at the book website), generate a ranking using Galago, and then calculate average precision, NDCG at 5 and 10, precision at 10, and the reciprocal rank by hand.

## 1.2 Approach

Galago version 3.10 was first downloaded from the Project Lemur Source Forge website, which can be found at the following URL: <https://sourceforge.net/projects/lemur/files/lemur/galago-3.10/>. The CACM document corpus was downloaded from the textbook's website, found here: <http://www.search-engines-book.com/collections/>. Galago was used to create an index of the CACM corpus and to run as a server to respond to queries on that index.

The `getrel.py` and `q83.py` scripts (found in Listings 3 and 4, respectively) was created to issue queries to the Galago search server using the Python Requests library [1]. The HTML responses were then parsed using the Python BeautifulSoup library [2], where the CACM document identifiers were extracted for use in calculating the different evaluation scores for the Galago ranking.

The query used was from the CACM query set, number 10, and only the first 1000 retrieved documents were considered when calculating all scores for this experiment.

### 1.2.1 Initial Precision and Recall Calculations

Precision and Recall were calculated with the following equations:

$$Recall = \frac{|A \cap B|}{|A|}$$
$$Precision = \frac{|A \cap B|}{|B|}$$

In these equations,  $A$  is the relevant set of documents for the query, and  $B$  is the set of retrieved documents.

### 1.2.2 Calculating Precision at Specific Rankings

A list of precision values was created by calculating the cumulative precision at each document ranking with the set of retrieved documents up to that ranking.

### 1.2.3 Calculating Average Precision

Average precision was calculated by adding the precision at each retrieval ranking position for documents which are part of  $A \cap B$ , or the set of retrieved documents that are relevant, and then dividing by the size of that set to obtain the average. This can also be described as the area under the precision-recall curve, which can be expressed as the following summation:

$$AveP = \sum_{k=1}^n P(k) \Delta r(k)$$

where  $k$  is the rank in the sequence of retrieved documents,  $n$  is the number of retrieved documents,  $P(k)$  is the precision at cut-off  $k$  in the list, and  $\Delta r(k)$  is the change in recall from items  $k - 1$  to  $k$ .

### 1.2.4 Calculating Normalized Discounted Cumulative Gain (NDCG)

First, discounted cumulative gain at rank  $p$  ( $DCG_p$ ) was calculated with the following formula:

$$DCG_p = rel_1 + \sum_{i=2}^p \frac{rel_i}{\log_2 i}$$

The ideal discounted cumulative gain at rank  $p$  ( $IDCG_p$ ) is a simple series, expressed as:

$$IDCG_p = 1 + \sum_{i=2}^p \frac{1}{\log_2 i}$$

Finally, normalized discounted cumulative gain at rank  $p$  ( $NDCG_p$ ) is expressed as:

$$NDCG_p = \frac{DCG_p}{IDCG_p}$$

with  $rel_i$  being the relevancy for document  $i$  in the retrieval ranking. For this experiment, this value is either 0 or 1.

### 1.2.5 Calculating Reciprocal Rank

Reciprocal rank is defined as the reciprocal of the rank at which the first relevant document is found, so if the  $3^{rd}$  document in the retrieval ranking list is the first relevant document, the reciprocal rank is  $\frac{1}{3}$ .

## 1.3 Results

After building the index, CACM query 10 was processed by the `getrel.py` script, the output of which can be found in Listing 1. This script calculates all the values shown in Table 1, which are all of the required values for the question.

```
1 [mchaney@mchaney-l getrel]$ python q8.3.py -n 1000 -q 10
2 query 10
3 query: parallel languages languages for parallel computation
4 precision: 0.027
5 recall: 0.771428571429
6 precision @10: 0.9
7 NDCG @5: 1.0
8 NDCG @10: 0.942709999032
9 avg precision: 0.697677898817
10 reciprocal rank: 1.0
```

Listing 1: Output from running the `getrel.py` script for queries 1 and 10 from the CACM collection.

Query #	Avg. Prec.	NDCG @5	NDCG @10	Prec. @10	Recip. Rank
10	0.697677898817	1.0	0.942709999032	0.9	1.0

Table 1: Calculations for CACM query 10 from top 1000 retrieved documents.

## 2 Question 8.4

### 2.1 Question

For two queries in the CACM collection, generate two uninterpolated recall-precision graphs, a table of interpolated precision values at standard recall levels, and the average interpolated recall-precision graph.

### 2.2 Approach

Using the `getrel.py`, `q84.py` and `graphs.R` scripts, found in Listings 3, 5 and 6.

```
59 def run(rel, retri, func):  
60     rr = []  
61     for i in range(1, len(retri)+1):  
62         rr.append(func(rel, retri[:i]))  
63     return rr
```

Listing 2: The run function

#### 2.2.1 Generating the Uninterpolated Recall-Precision Graph

The generated graph is shown in Figure ??.

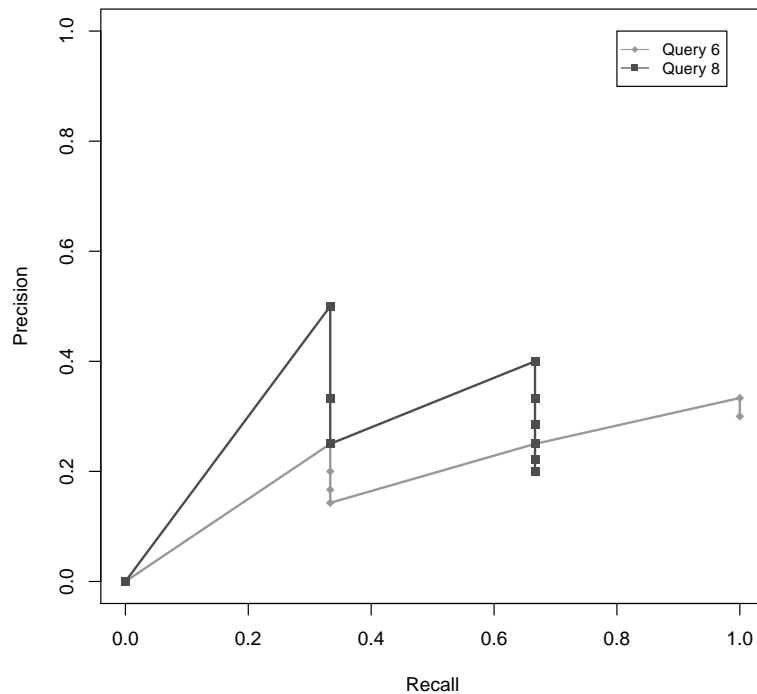


Figure 1: Uninterpolated Recall-Precision Graph for CACM Queries 6 and 8.

### 2.2.2 Generating the Table of Interpolated Precision Values

### 2.2.3 Generating the Average Interpolated Recall-Precision Graph

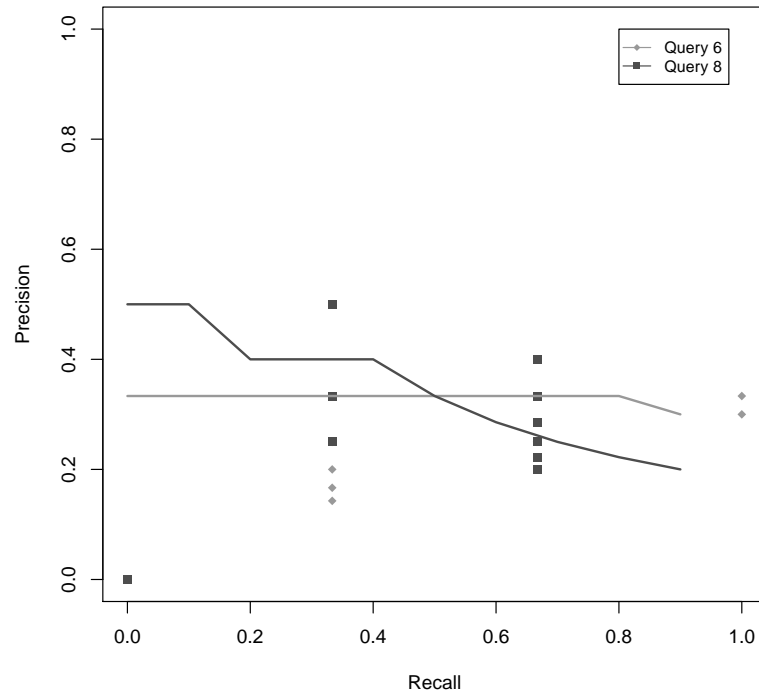


Figure 2: Average Interpolated Recall-Precision Graph for CACM Queries 6 and 8.

## 3 Appendix

### 3.1 Code listings

```
1 import argparse
2 import re
3 import requests
4 import xmltodict
5 import numpy as np
6 from math import log
7 from bs4 import BeautifulSoup
8
9
10 def parseargs():
11     parser = argparse.ArgumentParser()
12     parser.add_argument('-p', '--port', type=int, default=42247, help='the galago search
        server port')
13     parser.add_argument('-q', '--qnum', type=int, default=10, help='the query number to use')
14     parser.add_argument('-n', type=int, default=10, help='the number of results to retrieve')
15     return parser.parse_args()
16
17 args = parseargs()
18
19 def buildrel():
20     rel = {}
21     for line in open('cacm.rel').readlines():
22         q, _, doc, _ = line.split()
23         if q not in rel:
24             rel[q] = []
25         rel[q].append(int(doc.split('-')[1]))
26     return rel
27
28 def buildqueries():
29     with open('cacm.query.xml') as fd:
30         return xmltodict.parse(fd.read())
31
32 REL = buildrel()
33 QUERIES = buildqueries()
34 RE = re.compile('/home/mchaney/workspace/edu/cs834-f16/assignments/assignment4/code/cacm/
        docs/CACM-([d]+).html')
35 ID = {'id': 'result'}
36 URL = 'http://0.0.0.0:{0}/search'
37 QUERY1 = 'what articles exist which deal with tss time sharing system an operating system
        for ibm computers'
38 PDICT = {'q': QUERY1, 'start': 0, 'n': args.n}
39
40 def query(qstr, port=args.port):
41     PDICT['q'] = qstr
42     PDICT['n'] = args.n
43     res = requests.get(URL.format(port), params=PDICT)
44     if not res.ok:
45         return None
46     soup = BeautifulSoup(res.text, 'html.parser')
47     return [int(RE.match(href.text).groups()[0]) for href in soup.select("#result a")]
48
49 def recall(rel, retr):
50     relset = set(rel)
51     retrset = set(retr)
52     return float(len(relset.intersection(retrset))) / len(relset)
53
54 def precision(rel, retr):
55     relset = set(rel)
56     retrset = set(retr)
57     return float(len(relset.intersection(retrset))) / len(retrset)
58
59 def run(rel, retr, func):
60     rr = []
61     for i in range(1, len(retr)+1):
62         rr.append(func(rel, retr[:i]))
63     return rr
64
65 def avg(rel, retr, func):
66     prun = run(rel, retr, func)
```

```

67     res = []
68     for i in range(len(retr)):
69         if retr[i] in rel:
70             res.append(prun[i])
71     return float(sum(res))/len(res)
72
73 def getrel(rel, retr, i):
74     return 1 if retr[i] in rel else 0
75
76 def DCG(rel, retr, p):
77     sum = 0
78     for i in range(2, p+1):
79         sum += float(getrel(rel, retr, i-1)) / log(i, 2)
80     return getrel(rel, retr, 0) + sum
81
82 def IDCG(p):
83     sum = 0
84     for i in range(2, p+1):
85         sum += 1 / log(i, 2)
86     return 1 + sum
87
88 def NDCG(rel, retr, p):
89     dcg = DCG(rel, retr, p)
90     idcg = IDCG(p)
91     return dcg / idcg
92
93 def reciprank(rel, retr):
94     for i in range(1, len(retr)+1):
95         if retr[i-1] in rel:
96             return 1.0 / i
97     return 0.0
98
99 def ipr(rrun, prun):
100     res = []
101     res.append(max(prun))
102     for i in np.arange(0.1, 1, 0.1):
103         res.append(max(prun[int(i*10):]))
104     return np.arange(0, 1, 0.1).tolist(), res
105
106 def getquery(qnum):
107     return QUERIES['parameters'][ 'query' ][qnum-1][ 'text' ]
108
109 def process(qnum):
110     qstr = getquery(qnum)
111     retr = query(qstr)
112     rel = REL[str(qnum)]
113     prun = run(rel, retr, precision)
114     rrun = run(rel, retr, recall)
115     prec = precision(rel, retr)
116     rec = recall(rel, retr)
117     avgprec = avg(rel, retr, precision)
118     ndcg5 = NDCG(rel, retr, 5)
119     ndcg10 = NDCG(rel, retr, 10)
120     recip = reciprank(rel, retr)
121     return qnum, qstr, retr, rel, prun, rrun, prec, rec, ndcg5, ndcg10, avgprec, recip
122
123 def printresults(qnum, qstr, retr, rel, prun, rrun, prec, rec, ndcg5, ndcg10, avgprec, recip):
124     print 'query {0}'.format(qnum)
125     print 'query: {0}'.format(qstr)
126     if args.n == 10:
127         print 'relevant: {0}'.format(rel)
128         print 'retrieved: {0}'.format(retr)
129         print 'p-run: {0}'.format(prun)
130         print 'r-run: {0}'.format(rrun)
131         print 'precision: {0}'.format(prec)
132         print 'recall: {0}'.format(rec)
133         print 'precision @10: {0}'.format(prun[9])
134         print 'NDCG @5: {0}'.format(ndcg5)
135         print 'NDCG @10: {0}'.format(ndcg10)
136         print 'avg precision: {0}'.format(avgprec)
137         print 'reciprocal rank: {0}'.format(recip)

```

Listing 3: getrel.py



```

1 from getrel import *
2
3 TABLE = """\begin{table}[h!]
4 \centering
5 \begin{tabular}{c | c | c | c | c | c | c | c}
6 \hline
7 Query \# & Avg. Prec. & NDCG @5 & NDCG @10 & Prec. @10 & Recip. Rank & \\\
8 \hline
9 {0} & {1} & {2} & {3} & {4} & {5} & \\\
10 \hline
11 \end{tabular}
12 \caption{Calculations for CACM query {6} from top {7} retrieved documents.}
13 \label{tab:query20}
14 \end{table}
15 """
16
17 def printtab(qnum, qstr, retr, rel, prun, rrun, prec, rec, ndcg5, ndcg10, avgprec, recip):
18     fname = 'query{0}.tab'.format(qnum)
19     with open(fname, 'w') as fd:
20         fd.write(TABLE.format(qnum, avgprec, ndcg5, ndcg10, prun[9], recip, qnum, args.n))
21
22 results = process(args.qnum)
23 printresults(*results)
24 printtab(*results)

```

Listing 4: q83.py

```

1 from getrel import *
2
3 def printdata(rrun, prun, fname):
4     with open(fname, 'w') as fd:
5         zipped = zip(rrun, prun)
6         for z in zipped:
7             fd.write('{0}\t{1}\n'.format(z[0], z[1]))
8
9 results = process(args.qnum)
10 printresults(*results)
11 qnum, qstr, retr, rel, prun, rrun, prec, rec, ndcg5, ndcg10, avgprec, recip = results
12 printdata(rrun, prun, 'urpg{0}.dat'.format(qnum))
13
14 irrun, iprun = ipr(rrun, prun)
15 printdata(irrun, iprun, 'ipr{0}.dat'.format(qnum))

```

Listing 5: q84.py

```

1 plotone <- function(data, qnum) {
2   pdf(paste('urpg', qnum, '.pdf', sep=''))
3   plot(data, type='o', pch=15, ylim=c(0,1), xlim=c(0,1),
4         main=paste("Recall-Precision Graph for CACM Query ", qnum, sep=""),
5         ylab="Precision", xlab="Recall")
6   dev.off()
7 }
8 urpgraph <- function(d1, d2, fname) {
9   pdf(fname)
10  plot(d1, lwd=2, type='o', pch=18, ylim=c(0,1), xlim=c(0,1), col="gray60",
11        ylab="Precision", xlab="Recall")
12  lines(d2, lwd=2, type="o", pch=15, col="gray30")
13  legend(0.8, 1, c('Query 6', 'Query 8'), cex=0.8,
14        col=c('gray60', 'gray30'), lty=c(1,1), pch=c(18,15))
15  dev.off()
16 }
17 iprgraph <- function(d1, d2, id1, id2, fname) {
18  pdf(fname)
19  plot(d1, lwd=2, type="p", pch=18, ylim=c(0,1), xlim=c(0,1), col="gray60",
20        ylab="Precision", xlab="Recall")
21  lines(d2, lwd=2, type="p", pch=15, col="gray30")
22  lines(id1, lwd=2, type="l", col="gray60")
23  lines(id2, lwd=2, type="l", col="gray30")
24  legend(0.8, 1, c('Query 6', 'Query 8'), cex=0.8,
25        col=c('gray60', 'gray30'), lty=c(1,1), pch=c(18,15))
26  dev.off()
27 }
28 args = commandArgs(trailingOnly=TRUE)
29
30 d1 <- read.table(paste('urpg', args[1], '.dat', sep=''))
31 d2 <- read.table(paste('urpg', args[2], '.dat', sep=''))
32 plotone(d1, args[1])
33 plotone(d2, args[2])
34 urpgraph(d1, d2, paste('urpg', args[1], '', args[2], '.pdf', sep=''))
35
36 id1 <- read.table(paste('ipr', args[1], '.dat', sep=''))
37 id2 <- read.table(paste('ipr', args[2], '.dat', sep=''))
38 iprgraph(d1, d2, id1, id2, paste('ipr', args[1], '', args[2], '.pdf', sep=''))

```

Listing 6: Script used to generate the recall-precision graphs

## 4 References

- [1] Kenneth Reitz. Requests: HTTP for Humans. Available at <http://docs.python-requests.org/en/master/>. Accessed: 2016/09/20.
- [2] Leonard Richardson. Beautiful Soup. Available at: <https://www.crummy.com/software/beautifulsoup/>. Accessed: 2016/09/20.