# Assignment 4

## Fall 2016
## CS834 Introduction to Information Retrieval
## Dr. Michael Nelson

Mathew Chaney

December 4, 2016

# Contents

# Listings

# List of Figures

# List of Tables

# 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 2 and 3, 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 Beautiful Soup 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 2, 4 and 5 were created to complete this task.

## 2.3 Results

### 2.3.1 Uninterpolated Recall-Precision Graph

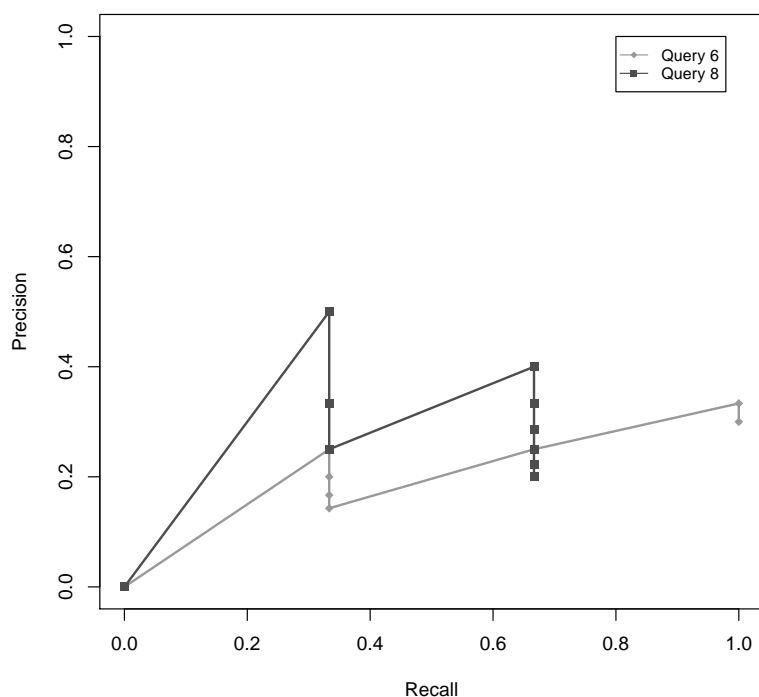The uninterpolated recall-precision graph is shown in Figure 1.



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

### 2.3.2 Interpolated Precision

The graph for the interpolated precision at standard recall values is shown in Figure 2 and the table of the values for each query, including the averages, is shown in Table 2.
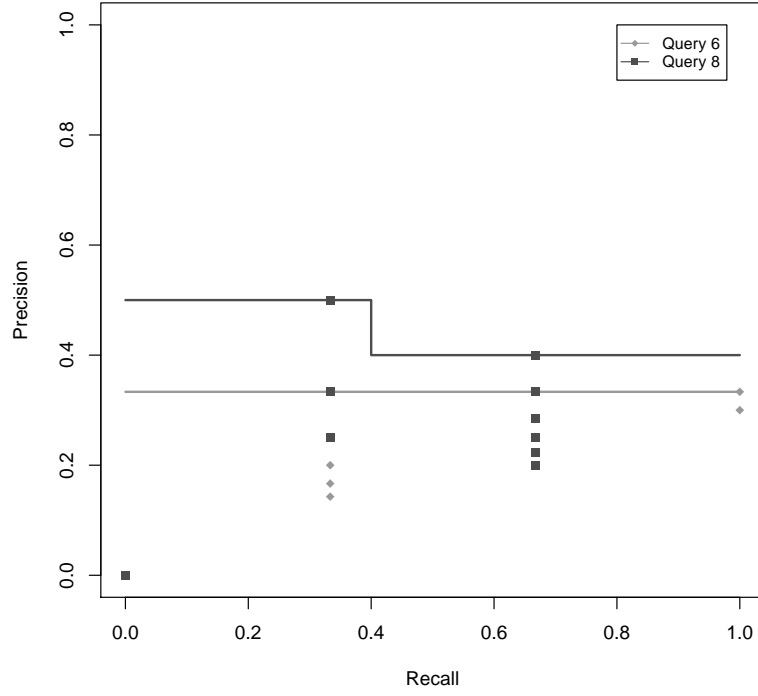
Figure 2: Graph of interpolated precision at standard recall values for queries 6 and 8.

| Recall | 0.0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1.0 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Query 6 | 0.333 | 0.333 | 0.333 | 0.333 | 0.333 | 0.333 | 0.333 | 0.333 | 0.333 | 0.333 | 0.333 |
| Query 8 | 0.5 | 0.5 | 0.5 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 |
| Average | 0.417 | 0.417 | 0.417 | 0.417 | 0.367 | 0.367 | 0.367 | 0.367 | 0.367 | 0.367 | 0.367 |

Table 2: Interpolated precision at standard recall values.

### 2.3.3 Average Interpolated Precision

The graph of the average interpolated precision at standard recall values for queries 6 and 8 can be found in Figure 3.
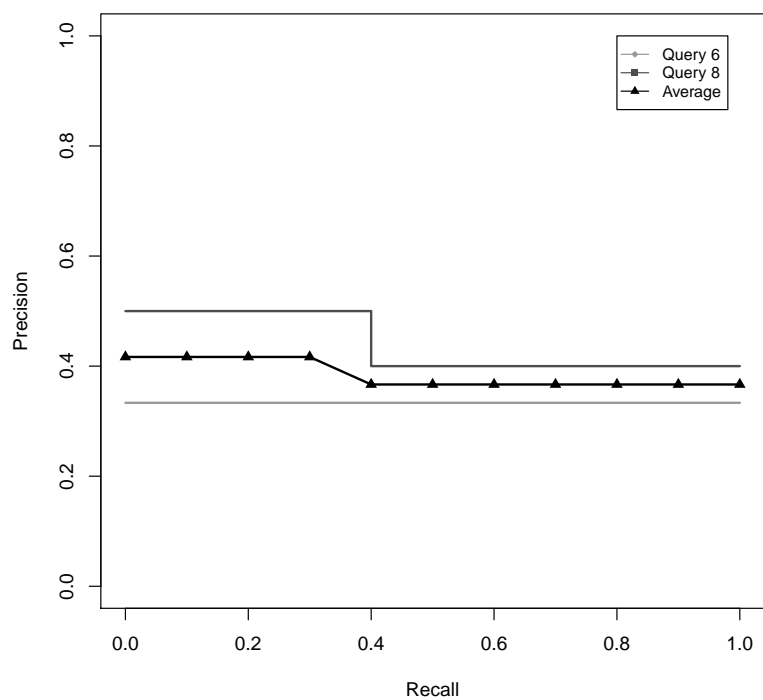
Figure 3: Average interpolated recall-precision graph for CACM Queries 6 and 8.

# 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='galago server port')
13     parser.add_argument('-q', '--qnum', nargs='+', type=int, default=[6, 8], help='query
           number')
14     parser.add_argument('-n', type=int, default=10, help='number of retrieval pages')
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, precision)
67     res = []
68     for i in range(len(retr)):
```

```python
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     for i in np.arange(0, 1.1, .1):
102         for j in range(len(rrun)):
103             if rrun[j] > i:
104                 idx = j
105                 break
106         res.append(max(prun[idx:]))
107     return np.arange(0, 1.1, 0.1), res
108
109 def getquery(qnum):
110     return QUERIES['parameters']['query'][qnum-1]['text']
111
112 def process(qnum):
113     qstr = getquery(qnum)
114     retr = query(qstr)
115     rel = REL[str(qnum)]
116     prun = run(rel, retr, precision)
117     rrun = run(rel, retr, recall)
118     prec = precision(rel, retr)
119     rec = recall(rel, retr)
120     avgprec = avg(rel, retr, precision)
121     ndcg5 = NDCG(rel, retr, 5)
122     ndcg10 = NDCG(rel, retr, 10)
123     recip = reciprank(rel, retr)
124     return qnum, qstr, retr, rel, prun, rrun, prec, rec, ndcg5, ndcg10, avgprec, recip
125
126 def printresults(qnum, qstr, retr, rel, prun, rrun, prec, rec, ndcg5, ndcg10, avgprec, recip
        ):
127     print 'query {0}'.format(qnum)
128     print 'query: {0}'.format(qstr)
129     if args.n == 10:
130         print 'relevant: {0}'.format(rel)
131         print 'retrieved: {0}'.format(retr)
132         print 'p-run: {0}'.format(prun)
133         print 'r-run: {0}'.format(rrun)
134     print 'precision: {0}'.format(prec)
135     print 'recall: {0}'.format(rec)
136     print 'precision @10: {0}'.format(prun[9])
137     print 'NDCG @5: {0}'.format(ndcg5)
138     print 'NDCG @10: {0}'.format(ndcg10)
139     print 'avg precision: {0}'.format(avgprec)
140     print 'reciprocal rank: {0}'.format(recip)
```

Listing 2: getrel.py

```python
from getrel import *

TABLE = """\\begin{{table}}[h!]
\\centering
\\begin{{tabular}}{{ | c | c | c | c | c | c | }}
\\hline
Query \# & Avg. Prec. & NDCG @5 & NDCG @10 & Prec. @10 & Recip. Rank \\\\
\\hline
{0} & {1} & {2} & {3} & {4} & {5} \\\\
\\hline
\\end{{tabular}}
\\caption{{Calculations for CACM query {6} from top {7} retrieved documents.}}
\\label{{tab:query20}}
\\end{{table}}
"""

def printtab(qnum, qstr, retr, rel, prun, rrun, prec, rec, ndcg5, ndcg10, avgprec, recip):
    fname = 'query{0}.tab'.format(qnum)
    with open(fname, 'w') as fd:
        fd.write(TABLE.format(qnum, avgprec, ndcg5, ndcg10, prun[9], recip, qnum, args.n))

for qnum in args.qnum:
    results = process(qnum)
    printresults(*results)
    printtab(*results)
```

Listing 3: q83.py

```python
from getrel import *

def printdata(rrun, prun, fname):
    with open(fname, 'w') as fd:
        zipped = zip(rrun, prun)
        for z in zipped:
            fd.write('{0}\t{1}\n'.format(z[0], z[1]))

HEAD = """\\begin{table}[H]
\\centering
\\begin{tabular}{ l l l l l l l l l l l l l }
Recall & 0.0 & 0.1 & 0.2 & 0.3 & 0.4 & 0.5 & 0.6 & 0.7 & 0.8 & 0.9 & 1.0 \\\\
\\cline{2-12}
"""

ROW = """{0} & {1:.3g} & {2:.3g} & {3:.3g} & {4:.3g} & {5:.3g} & {6:.3g} & {7:.3g} & {8:.3g}
    & {9:.3g} & {10:.3g} & {11:.3g} \\\\
\\cline{{2-12}}
"""

TAIL = """\\end{tabular}
\\caption{.}
\\label{tab:ipr68}
\\end{table}
"""

def printtable(iprl):
    with open('iptab.tex', 'w') as fd:
        fd.write(HEAD)
        for iprun, qnum in iprl:
            fd.write(ROW.format('Query {0}'.format(qnum), *iprun))
        avg = [float(sum(col))/len(col) for col in zip(*[col[0] for col in iprl])]
        printdata(np.arange(0, 1.1, .1), avg, 'avg.dat')
        fd.write(ROW.format('Average', *avg))
        fd.write(TAIL)

iprl = []
for qnum in args.qnum:
    results = process(qnum)
    printresults(*results)
    qnum, qstr, retr, rel, prun, rrun, prec, rec, ndcg5, ndcg10, avgprec, recip = results
    printdata(rrun, prun, 'urpg{0}.dat'.format(qnum))
    irrun, iprun = ipr(rrun, prun)
    printdata(irrun, iprun, 'ipr{0}.dat'.format(qnum))
    iprl.append((iprun, qnum))
printtable(iprl)
```

```
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(id1, lwd=2, type="s", col="gray60")
22     lines(d2, lwd=2, type="p", pch=15, col="gray30")
23     lines(id2, lwd=2, type="s", 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 aipgraph <- function(avg, id1, id2, fname) {
29     pdf(fname)
30     plot(avg, lwd=2, type="l", ylim=c(0,1), xlim=c(0,1), col="black",
31         ylab="Precision", xlab="Recall")
32     lines(avg, lwd=2, type="p", pch=17, col="black")
33     lines(id1, lwd=2, type="s", col="gray60")
34     lines(id2, lwd=2, type="s", col="gray30")
35     legend(0.8, 1, c('Query 6', 'Query 8', 'Average'), cex=0.8,
36         col=c('gray60', 'gray30', 'black'), lty=c(1,1,1), pch=c(18,15,17))
37     dev.off()
38 }
39
40 args = commandArgs(trailingOnly=TRUE)
41
42 d1 <- read.table(paste('urpg', args[1], '.dat', sep=''))
43 d2 <- read.table(paste('urpg', args[2], '.dat', sep=''))
44 plotone(d1, args[1])
45 plotone(d2, args[2])
46 urpgraph(d1, d2, paste('urpg', args[1], '', args[2], '.pdf', sep=''))
47
48 id1 <- read.table(paste('ipr', args[1], '.dat', sep=''))
49 id2 <- read.table(paste('ipr', args[2], '.dat', sep=''))
50 iprgraph(d1, d2, id1, id2, paste('ipr', args[1], '', args[2], '.pdf', sep=''))
51
52 avg <- read.table('avg.dat')
53 aipgraph(avg, id1, id2, paste('aipr', args[1], args[2], '.pdf', sep=''))
```

Listing 5: 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.