

Assignment 4

Fall 2016

CS834 Introduction to Information Retrieval

Dr. Michael Nelson

Mathew Chaney

December 3, 2016

Contents

1	Question 8.3	3
1.1	Question	3
1.2	Approach	3
1.3	Results	4
2	Question 8.4	5
2.1	Question	5
2.2	Approach	5
2.3	Results	5
3	Appendix	6
3.1	Code listings	6
4	References	9

List of Figures

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 2 and 5, 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

2.2.1 Generating the Uninterpolated Recall-Precision Graph

2.2.2 Generating the Table of Interpolated Precision Values

2.2.3 Generating the Average Interpolated Recall-Precision Graph

2.3 Results

3 Appendix

3.1 Code listings

```
1 import argparse
2 import re
3 import requests
4 import xmltodict
5 from math import log
6 from bs4 import BeautifulSoup
7
8
9 def parseargs():
10     parser = argparse.ArgumentParser()
11     parser.add_argument('-q', '--qnum', type=int, default=10, help='the query number to use')
12     parser.add_argument('-n', type=int, default=10, help='the number of results to retrieve')
13     return parser.parse_args()
14
15 args = parseargs()
16
17 def buildrel():
18     rel = {}
19     for line in open('cacm.rel').readlines():
20         q, _, doc, _ = line.split()
21         if q not in rel:
22             rel[q] = []
23         rel[q].append(int(doc.split('-')[1]))
24     return rel
25
26 def buildqueries():
27     with open('cacm.query.xml') as fd:
28         return xmltodict.parse(fd.read())
29
30 REL = buildrel()
31 QUERIES = buildqueries()
32 RE = re.compile('/home/mchaney/workspace/edu/cs834-f16/assignments/assignment4/code/cacm/docs/CACM-([\\d]+).html')
33 ID = {'id': 'result'}
34 URL = 'http://0.0.0.0:{0}/search'
35 QUERY1 = 'what articles exist which deal with tss time sharing system an operating system for ibm computers'
36 PDICT = {'q': QUERY1, 'start': 0, 'n': args.n}
37
38 def query(qstr, port=54312):
39     PDICT['q'] = qstr
40     PDICT['n'] = args.n
41     res = requests.get(URL.format(port), params=PDICT)
42     if not res.ok:
43         return None
44     soup = BeautifulSoup(res.text, 'html.parser')
45     return [int(RE.match(href.text).groups()[0]) for href in soup.select("#result a")]
46
47 def recall(rel, retr):
48     relset = set(rel)
49     retrset = set(retr)
50     return float(len(relset.intersection(retrset))) / len(relset)
51
52 def precision(rel, retr):
53     relset = set(rel)
54     retrset = set(retr)
55     return float(len(relset.intersection(retrset))) / len(retrset)
56
57 def run(rel, retr, func):
58     rr = []
59     for i in range(1, len(retr)+1):
60         rr.append(func(rel, retr[:i]))
61     return rr
62
63 def avg(rel, retr, func):
64     prun = run(rel, retr, func)
65     res = []
66     for i in range(len(retr)):
67         if retr[i] in rel:
```

```

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

```

Listing 2: getrel.py

```

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

```

Listing 3: q83.py

```

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

```

Listing 4: q84.py

```

1 args = commandArgs(trailingOnly=TRUE)
2 data <- read.table(paste('urpg', args[1], '.dat', sep=''))
3
4 ploturpg <- function(data) {
5     pdf(paste('urpg.pdf', sep=''))
6     plot(data, type='o', ylim=c(0,1), xlim=c(0,1), pch=15)
7     dev.off()
8 }
9
10 ploturpg(data)

```

Listing 5: Script used to generate the uninterpolated recall-precision graph

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.