

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	Approach	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 script (found in Listing 2) 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 getrel.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 Approach

Git er done

3 Appendix

3.1 Code listings

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

```

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

```

```

145     print 'avg precision: {0}'.format(avgprec)
146     print 'reciprocal rank: {0}'.format( recip)
147
148
149 TABLE = """\begin{table}[h!]
150 \centering
151 \begin{tabular}{| c | c | c | c | c | c | }
152 \hline
153 Query \# & Avg. Prec. & NDCG @5 & NDCG @10 & Prec. @10 & Recip. Rank \\\
154 \hline
155 {0} & {1} & {2} & {3} & {4} & {5} \\\
156 \hline
157 \end{tabular}
158 \caption{Calculations for CACM query {6} from top {7} retrieved documents.}
159 \label{tab:query20}
160 \end{table}
161 """
162
163 def printtab(qnum, qstr, retr, rel, prun, prec, rec, ndcg5, ndcg10, avgprec, recip):
164     fname = 'query{0}.tab'.format(qnum)
165     with open(fname, 'w') as fd:
166         fd.write(TABLE.format(qnum, avgprec, ndcg5, ndcg10, prun[9], recip, qnum, args.n))
167
168
169 results = process(args.qnum)
170 printresults(*results)
171 printtab(*results)

```

Listing 2: getrel.py

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.