

# **Assignment 4**

**Fall 2016**

**CS834 Introduction to Information Retrieval**

**Dr. Michael Nelson**

Mathew Chaney

December 3, 2016

## Contents

<b>1</b>	<b>Question 8.3</b>	<b>3</b>
1.1	Question . . . . .	3
1.2	Approach . . . . .	3
<b>2</b>	<b>Appendix</b>	<b>4</b>
2.1	Code listings . . . . .	4
<b>3</b>	<b>References</b>	<b>6</b>

## 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/>.

An index of the CACM corpus, downloaded from the book website <http://www.search-engines-book.com/collections/>, was created with Galago with the following command:

```
galago build -indexPath=cacm.index -inputPath=docs -server=true
```

This index was used with Galago running as a search engine with the following command:

```
galago search -index=cacm.index
```

The getrel.py script was created to issue queries to the running Galago search server using the Python Requests library [1]. The HTML responses were then parsed using the Python BeautifulSoup library [?], where the CACM document identifiers were extracted for use in calculating the different evaluation scores for the Galago ranking.

Precision and Recall were calculated with the following equations:

$$Recall = \frac{|A \cap B|}{|A|} \quad (1)$$

$$Precision = \frac{|A \cap B|}{|B|} \quad (2)$$

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

### 1.2.1 Calculating Average Precision for CACM Query 1

Query 1 and 10 were used for this question. Average precision was calculated

## 2 Appendix

### 2.1 Code listings

```
1 #!/usr/bin/python
2
3 import re
4 import requests
5 from bs4 import BeautifulSoup
6 from pprint import pprint as pp
7
8
9 def buildrel():
10     rel = {}
11     for line in open('cacm.rel').readlines():
12         q, _, doc, _ = line.split()
13         if q not in rel:
14             rel[q] = []
15         rel[q].append(doc)
16     return rel
17
18
19 RE = re.compile('/home/mchaney/workspace/edu/cs834-f16/assignments/assignment4/code/cacm/
20 docs/(CACM-[\\d]+).html')
21 ID = {'id': 'result'}
22 URL = 'http://0.0.0.0:{0}/search'
23 QUERY1 = 'what articles exist which deal with tss time sharing system an operating system
24 for ibm computers'
25 QUERY10 = 'parallel languages languages for parallel computation'
26 PDICT = {'q': QUERY1}
27
28 def query(query, port=54312):
29     PDICT['q'] = query
30     res = requests.get(URL.format(port), params=PDICT)
31     if not res.ok:
32         return None
33     soup = BeautifulSoup(res.text, 'html.parser')
34     return [RE.match(href.text).groups()[0] for href in soup.select("#result a")]
35
36
37 def recall(rel, retr):
38     A = set(rel)
39     B = set(retr)
40     return float(len(A.intersection(B))) / len(A)
41
42
43 def precision(rel, retr):
44     A = set(rel)
45     B = set(retr)
46     return float(len(A.intersection(B))) / len(B)
47
48
49 def run(rel, retr, func):
50     rr = []
51     for i in range(1, len(retr)):
52         rr.append(func(rel, retr[:i]))
53     return rr
54
55
56 def avg(rel, retr, func):
57     rr = run(rel, retr, func)
58     return sum(rr)/len(rr)
59
60
61 rel = buildrel()
62
63 print 'query:', QUERY1
64 retr = query(QUERY1)
65 print 'retrieved:', retr
66 print 'relevant:', rel['1']
67 print 'precision:', precision(retr, rel['1'])
68 print 'recall:', recall(retr, rel['1'])
69 print 'average precision:', avg(retr, rel['1'], precision)
70 print
71 print 'query:', QUERY10
```

```

70| retr = query(QUERY10)
71| print 'retrieved:', retr
72| print 'relevant:', rel['10']
73| print 'precision:', precision(retr, rel['10'])
74| print 'recall:', recall(retr, rel['10'])
75| print 'average precision:', avg(retr, rel['10'], precision)

```

Listing 1: getrel.py

```

1 (skipped debug output)
2 Stage parsePostings completed with 0 errors.
3 Stage writeFields completed with 0 errors.
4 Stage writeNames completed with 0 errors.
5 Stage writeLengths completed with 0 errors.
6 Stage writeNamesRev completed with 0 errors.
7 Stage writeCorpusKeys completed with 0 errors.
8 Stage writePostings completed with 0 errors.
9 Stage writePostings-krovetz completed with 0 errors.
10 2016-12-02 07:30:20.267:INFO:oejs.ServerConnector:main: Stopped ServerConnector@3b4d8f47{
    HTTP/1.1,[http/1.1]}{0.0.0.0:50935}
11 Done Indexing.
12 - 0.00 Hours
13 - 0.18 Minutes
14 - 11.01 Seconds
15 Documents Indexed: 3204.

```

Listing 2: Output from building the CACM index.

```

1 [mchaney@mchaney-l getrel]$ python getrel.py
2 query: what articles exist which deal with tss time sharing system an operating system for
    ibm computers
3 retrieved: [u'CACM-1410', u'CACM-1827', u'CACM-0397', u'CACM-1280', u'CACM-2319', u'CACM-
    -1938', u'CACM-1908', u'CACM-1885', u'CACM-1071', u'CACM-2535']
4 relevant: ['CACM-1410', 'CACM-1572', 'CACM-1605', 'CACM-2020', 'CACM-2358']
5 precision: 0.2
6 recall: 0.1
7 average precision: 0.520833333333
8
9 query: parallel languages languages for parallel computation
10 retrieved: [u'CACM-2785', u'CACM-1811', u'CACM-1262', u'CACM-0950', u'CACM-2895', u'CACM-
    -2700', u'CACM-2851', u'CACM-1747', u'CACM-2289', u'CACM-2266']
11 relevant: ['CACM-46', 'CACM-141', 'CACM-392', 'CACM-950', 'CACM-1158', 'CACM-1198', 'CACM-
    -1262', 'CACM-1380', 'CACM-1471', 'CACM-1601', 'CACM-1613', 'CACM-1747', 'CACM-1795', '
    CACM-1811', 'CACM-2060', 'CACM-2150', 'CACM-2256', 'CACM-2289', 'CACM-2342', 'CACM-
    -2376', 'CACM-2433', 'CACM-2618', 'CACM-2664', 'CACM-2685', 'CACM-2700', 'CACM-2714', '
    CACM-2777', 'CACM-2785', 'CACM-2851', 'CACM-2895', 'CACM-2896', 'CACM-2912', 'CACM-
    -3039', 'CACM-3075', 'CACM-3156']
12 precision: 0.228571428571
13 recall: 0.8
14 average precision: 0.155878721569

```

Listing 3: Output from running the getrel.py script for queries 1 and 10 from the CACM collection.

### 3 References

- [1] Kenneth Reitz. Requests: HTTP for Humans. Available at <http://docs.python-requests.org/en/master/>. Accessed: 2016/09/20.