

Assignment 2

Fall 2016

CS834 Introduction to Information Retrieval

Dr. Michael Nelson

Mathew Chaney

October 11, 2016

Contents

1	Question 4.1	3
1.1	Question	3
1.2	Resources	3
1.3	Answer	3
2	Appendix A	5
3	References	8

List of Figures

1	Word Counts for Small Wikipedia Corpus	3
2	Bigram Counts for Small Wikipedia Corpus	4

Listings

1	wc.py	5
2	buildgraphs.R	7

List of Tables

1 Question 4.1

1.1 Question

Plot rank-frequency curves (using a log-log graph) for words and bigrams in the Wikipedia collection available through the book website (<http://www.search-engines-book.com>). Plot a curve for the combination of the two. What are the best values for the parameter c for each curve?

1.2 Resources

The textbook *Search Engines: Information Retrieval in Practice* [1], the Python programming language [2], the R programming language [?] and the BeautifulSoup python library [3] were used to answer this question.

1.3 Answer

The wc.py script 1 was used to locate each file of the Wikipedia collection obtained from the book download page, available at <http://www.search-engines-book.com>. The BeautifulSoup library was used to strip out the HTML tags and then the nltk library [?] was used to tokenize the text. The individual words were counted manually and the nltk library [?] was used to count the bigrams.

The word count graph can be found in Figure 1 and the bigram count graph can be found in Figure 2.

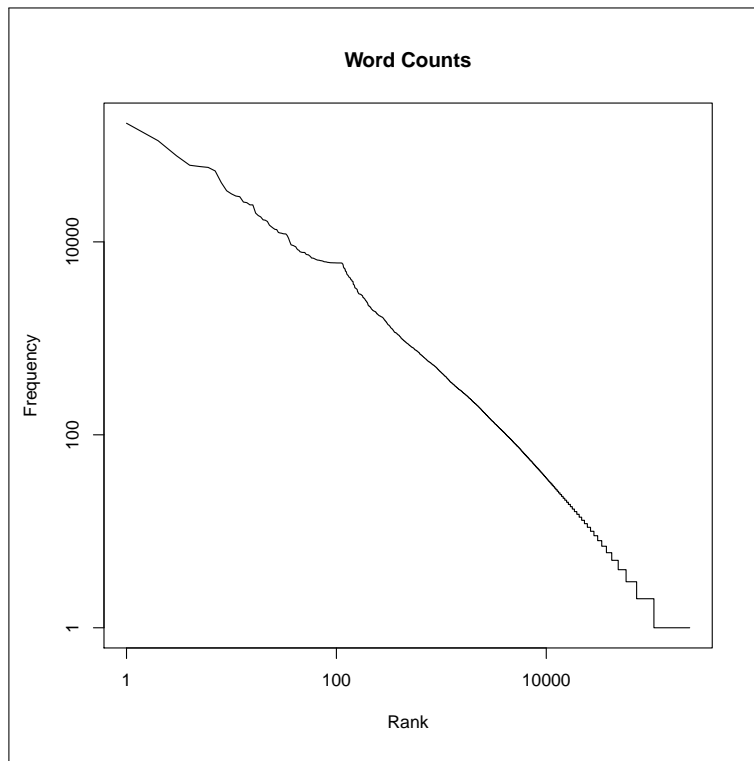


Figure 1: Word Counts for Small Wikipedia Corpus

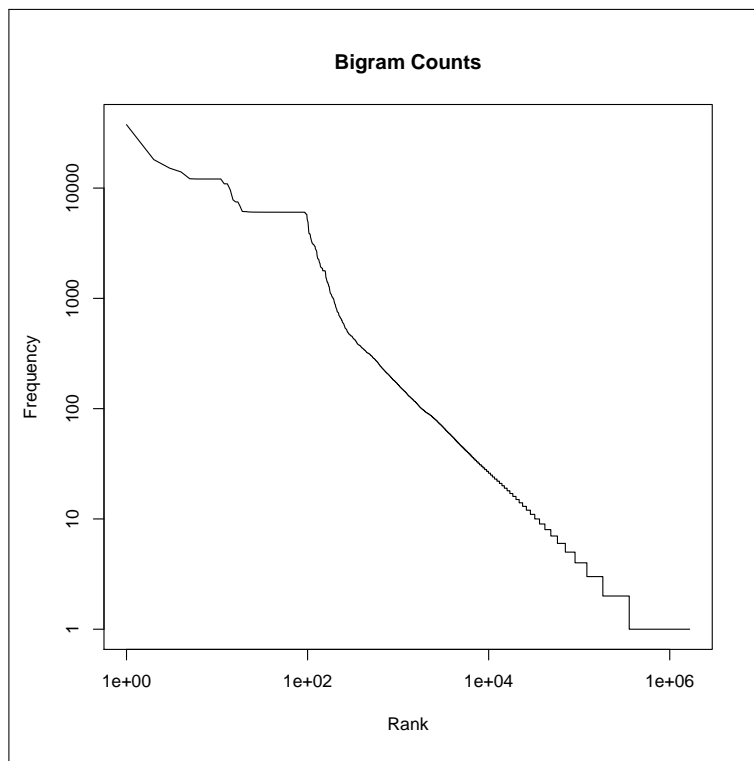


Figure 2: Bigram Counts for Small Wikipedia Corpus

2 Appendix A

```
1 #!/usr/bin/python
2
3 import argparse
4 import os
5 import operator
6 import sys
7 import nltk
8 from os.path import isdir, isfile
9 from bs4 import BeautifulSoup
10
11 class WordCounter(object):
12
13     def __init__(self, root):
14         self.tokenizer = nltk.RegexpTokenizer(r'\w+')
15         self.root = root
16         self.wmap = {}
17         self.bgmap = {}
18         self.filelist = []
19         self.visited = 0
20
21     def getfiles(self, folder=''):
22         items = os.listdir(self.root + folder)
23         for item in items:
24             filepath = self.root + folder + os.sep + item
25             if isfile(filepath):
26                 self.filelist.append(filepath)
27             elif isdir(filepath):
28                 self.getfiles(folder + os.sep + item)
29
30     def count(self, filepath):
31         sys.stdout.write("\rprocessing document #%i" % self.visited)
32         sys.stdout.flush()
33         with open(filepath) as infile:
34             text = BeautifulSoup(infile.read(), 'html.parser').get_text()
35             tokens = self.tokenizer.tokenize(text)
36             for s in tokens:
37                 if not self.wmap.has_key(s):
38                     self.wmap[s] = 0
39                 self.wmap[s] = self.wmap[s] + 1
40             for b in nltk.bigrams(tokens):
41                 if not self.bgmap.has_key(b):
42                     self.bgmap[b] = 0
43                 self.bgmap[b] = self.bgmap[b] + 1
44             self.visited = self.visited + 1
45
46     def writeresults(self):
47         with open('wordcount', 'w') as outfile:
48             for k, v in sorted(self.wmap.items(), key=operator.itemgetter(1), reverse=True):
49                 outfile.write(str(v) + '\t' + k.encode('utf-8') + '\n')
50         with open('bigramcount', 'w') as outfile:
51             for k, v in sorted(self.bgmap.items(), key=operator.itemgetter(1), reverse=True):
52                 outfile.write(str(v) + '\t' + k[0].encode('utf-8') + '\t' + k[1].encode('utf-8') + '\n')
53
54     def run(self):
55         print 'delving into "{0}"'.format(self.root)
56         self.getfiles()
57         print 'found {0} documents'.format(len(self.filelist))
58         map(self.count, self.filelist)
59         print '\nfound {0} words'.format(len(self.wmap))
60         print 'found {0} bigrams'.format(len(self.bgmap))
61         self.writesults()
62
63
64 if __name__ == '__main__':
65     parser = argparse.ArgumentParser('word count')
66     parser.add_argument('-root', '-r', help='the root directory for parsing', default='en')
67     args = parser.parse_args()
68
69     wc = WordCounter(args.root)
70     wc.run()
```

Listing 1: wc.py

```

1 #! /usr/bin/Rscript
2
3 plotgraph <- function(infile, outfile, title) {
4   data <- read.table(infile)
5   x <- seq(1, length(data$V1))
6   y <- data$V1
7
8   pdf(outfile)
9   plot(x, y, type='l', log='xy', main=title,
10        ylab='Frequency', xlab='Rank')
11   dev.off()
12 }
13
14 plotgraph('wordcount', 'wc.pdf', 'Word Counts')
15 plotgraph('bigramcount', 'bg.pdf', 'Bigram Counts')

```

Listing 2: buildgraphs.R

3 References

- [1] Bruce Croft, Donald Metzler, and Trevor Strohman. *Search Engines: Information Retrieval in Practice*. Pearson, first edition, February 2009.
- [2] The Python Programming Language. Available at: <https://www.python.org/>. Accessed: 2016/09/17.
- [3] Leonard Richardson. Beautiful Soup. Available at: <https://www.crummy.com/software/beautifulsoup/>. Accessed: 2016/09/20.