

Assignment 3

Fall 2016

CS834 Introduction to Information Retrieval

Dr. Michael Nelson

Mathew Chaney

November 9, 2016

Contents

1	Question 6.1	3
1.1	Question	3
1.2	Approach	3
2	Question 6.2	4
2.1	Question	4
2.2	Approach	4
2.3	Results	4
3	Question 6.3	5
3.1	Question	5
3.2	Approach	5
4	Question 6.5	6
4.1	Question	6
4.2	Answer	6
5	Question MLN2	7
5.1	Question	7
5.2	Approach	7
5.3	Results	7
6	Appendix	11
7	References	15

List of Figures

Listings

1	spelling.py example output	4
2	stem.py	11
3	data.py	12
4	spelling.py	12
5	calc.py	13

List of Tables

1	Calculated values for “running”	7
2	Calculated values for “calculation”	7
3	Calculated values for “color”	8
4	Calculated values for “horse”	8
5	Calculated values for “sky”	8
6	Calculated values for “railroad”	9
7	Calculated values for “calendar”	9
8	Calculated values for “airplane”	9
9	Calculated values for “ocean”	10
10	Calculated values for “bicycle”	10

1 Question 6.1

1.1 Question

Using the Wikipedia collection provided at the book website, create a sample of stem clusters by the following process:

1. Index the collection without stemming.
2. Identify the first 1,000 words (in alphabetical order) in the index.
3. Create stem classes by stemming these 1,000 words and recording which words become the same stem.
4. Compute association measures (Dice's coefficient) between all pairs of stems in each stem class. Compute co-occurrence at the document level.
5. Create stem clusters by thresholding the association measure. All terms that are still connected to each other form the clusters.

Compare the stem clusters to the stem classes in terms of size and the quality (in your opinion) of the groupings.

1.2 Approach

The `stem.py` script, found in Listing 2, was used to solve this problem.

2 Question 6.2

2.1 Question

Create a simple spelling corrector based on the noisy channel model. Use a single-word language model, and an error model where all errors with the same edit distance have the same probability. Only consider edit distances of 1 or 2. Implement your own edit distance calculator (example code can easily be found on the Web)

2.2 Approach

Peter Norvig's noisy channel spelling correction algorithm [1] was used as the basis for this solution. The `spelling.py` script, found in Listing 4, was created as an implementation of this algorithm. It was written with the Python programming language [2].

A large text file was downloaded from Mr. Norvig's website to calculate language model probability function $P(W)$. The words in the text file were counted and stored in a map that was compressed and saved on disk using the pickle python library [3].

$P(W)$ is calculated with the following formula:

$$P(W) = \frac{C_W}{N}$$

where C_W is the word count for word W and N is the sum of all word counts.

The process of determining a spelling correction is as follows:

1. Take the input word and determine all existing (correctly spelled) words with edit distance one and two.
2. With the assumption that shorter edit distances equate to a higher probability of being the correct intended word, select from the set of words from the previous step the one with the shortest edit distance and highest value for $P(W)$.

2.3 Results

Here is some sample output from the `spelling.py` script.

```
1 [mchaney@mchaney-1 spelling]$ ./spelling splling
2 selling
3 [mchaney@mchaney-1 spelling]$ ./spelling sweling
4 swelling
5 [mchaney@mchaney-1 spelling]$ ./spelling aacck
6 back
7 [mchaney@mchaney-1 spelling]$ ./spelling panaceu
8 palace
9 [mchaney@mchaney-1 spelling]$ ./spelling plaec
10 place
11 [mchaney@mchaney-1 spelling]$ ./spelling intrmdiate
12 intermediate
13 [mchaney@mchaney-1 spelling]$ ./spelling informatino
14 information
15 [mchaney@mchaney-1 spelling]$ ./spelling pretende
16 pretended
17 [mchaney@mchaney-1 spelling]$ ./spelling teh
18 the
```

Listing 1: `spelling.py` example output

3 Question 6.3

3.1 Question

Implement a simple pseudo-relevance feedback algorithm for the Galago search engine. Provide examples of the query expansions that your algorithm does, and summarize the problems and successes of your approach.

3.2 Approach

Here's a formula.

$$\frac{n_{ab}}{n_a + n_b}$$

4 Question 6.5

4.1 Question

Describe the snippet generation algorithm in Galago. Would this algorithm work well for pages with little text content? Describe in detail how you would modify the algorithm to improve it.

4.2 Answer

Snippet creation is done by the `SnippetGenerator` class. This class takes as parameters to its `getSnippet` method the document text as a `String` and a `Set` of `String` query terms, and returns a `String` that is a query-relevant snippet, or summary, of the document.

The snippet generator begins by turning the document text into a list of tokens for processing. The generator then parses these tokens, looking for query term matches, and when it finds a match, it creates a `SnippetRegion` object that stores the location within the document where the query term matched, plus five contextual terms preceding and following each term match. This equates to storing sentence fragments containing query terms.

After collecting all of the regions in the document containing a query term the generator begins constructing the final snippet by adding the `SnippetRegions` found from the previous step, combining those regions that overlap each other into larger regions, until a final list of `SnippetRegions` is created with total length in terms is no greater than $40 + \text{the length of the last } \text{SnippetRegion} \text{ added}$.

With the final list of `SnippetRegions` the algorithm builds an HTML string containing all the snippets concatenated together for rendering the snippet in a browser while adding `` tags around each query term match for emphasis.

This approach favors regions at the beginning of the document without regard to query context. One way to improve upon this method is to favor regions that contain more query terms. This can be done by counting the number of query terms found in the combined regions and then ordering the snippet generation based on the regions with the highest contained query term counts. This method could cut down the size of the final snippet by choosing regions that contain more query words within the normal extent of 5 terms per query word match, which would allow for a more concise summary of the website as it relates to the user query.

5 Question MLN2

5.1 Question

Using the small wikipedia example, choose 10 words and compute MIM, EMIM, chi square, dice association measures for full document & 5 word windows (cf. pp. 203-205)

5.2 Approach

The python script `calc.py`, found in Listing 5, was used to complete this task.

5.3 Results

Here is the output from running the `calc.py` script:

running			
<i>MIM</i>	<i>EMIM</i>	χ^2	<i>Dice</i>
Tootie	Tootie	long	ran
Mortes	Mortes	only	long
Mortem	Mortem	but	could
Alsab	Alsab	over	run
Titulus	Titulus	two	started
Cruguet	Cruguet	could	ever
defensed	defensed	had	changed
Vipiteno	Vipiteno	time	old
Velocisaurus	Velocisaurus	In	opening
Pedophilia	Pedophilia	into	end

Table 1: Calculated values for “running”

calculation			
<i>MIM</i>	<i>EMIM</i>	χ^2	<i>Dice</i>
unknot	unknot	proleptic	usefulness
Jabr	Jabr	Casull	Spoon
humbler	humbler	Exiguus	computed
Bcbell	Bcbell	usefulness	compute
Marxschen	Marxschen	Spoon	calculate
Ethiopic	Ethiopic	computed	formulas
reconciling	reconciling	falsify	proleptic
anthropologie	anthropologie	compute	Casull
dampens	dampens	calculate	Exiguus
provable	provable	formulas	falsify

Table 2: Calculated values for “calculation”

color			
<i>MIM</i>	<i>EMIM</i>	χ^2	<i>Dice</i>
roadrunners	roadrunners	Depreciated	Depreciated
Tootie	Tootie	param	param
SparrowsWing	SparrowsWing	Alter	red
equilateral	equilateral	Abilities	colors
Sleepwalking	Sleepwalking	ego	black
Editorials	Editorials	red	Comics
Alor	Alor	colors	infobox
Antaheen	Antaheen	white	white
mutantsHidden	mutantsHidden	black	image
Caucasoids	Caucasoids	NGV17	ego

Table 3: Calculated values for “color”

horse			
<i>MIM</i>	<i>EMIM</i>	χ^2	<i>Dice</i>
Alsab	Alsab	thoroughbred	Horse
Cruguet	Cruguet	Equestrianism	thoroughbred
haoma	haoma	Zafonic	Stakes
pompeux	pompeux	Stakes	Equestrianism
iro	iro	racehorse	Zafonic
Awaystay	Awaystay	racehorses	racehorse
Beaurepaire	Beaurepaire	Thoroughbred	racehorses
Jardim	Jardim	Horse	Thoroughbred
Agnihotra	Agnihotra	Harness	Trainer
Legate	Legate	Slipper	racing

Table 4: Calculated values for “horse”

sky			
<i>MIM</i>	<i>EMIM</i>	χ^2	<i>Dice</i>
mailings	mailings	binoculars	Astronomy
Hig	Hig	ChristalPalace	bright
Alor	Alor	calvus	wind
Jeremywn	Jeremywn	Arcus	items
Kert01	Kert01	incus	eclipse
Chikubasho	Chikubasho	mackerel	visible
Jabr	Jabr	æŨĜ	speeds
Sennen	Sennen	Achiu31	gravity
iro	iro	Colares	Telescope
Cucumber	Cucumber	Cycles	objects

Table 5: Calculated values for “sky”

railroad			
<i>MIM</i>	<i>EMIM</i>	χ^2	<i>Dice</i>
Timken	Timken	railroads	Railroad
Hegins	Hegins	Railroad	railroads
Sameerkale	Sameerkale	Railroads	Slambo
Contr��tle	Contr��tle	Slambo	rail
Friedensburg	Friedensburg	trackage	freight
WLVN	WLVN	freight	Railroads
Harrisonville	Harrisonville	rail	Railway
C420	C420	Railway	gauge
C425	C425	gauge	Lines
C424	C424	mae	train

Table 6: Calculated values for “railroad”

calendar			
<i>MIM</i>	<i>EMIM</i>	χ^2	<i>Dice</i>
27a	27a	Gregorian	Gregorian
S��renstam	S��renstam	liturgics	liturgical
Jabr	Jabr	Lunisolar	calendars
escalade	escalade	Tixity	lunar
Tankersley	Tankersley	Calendarists	Persia
Desinicization	Desinicization	commemorations	Dionysius
Kikadue	Kikadue	calendars	Calendar
Munaishy	Munaishy	liturgical	Frysk
Mandarina999	Mandarina999	Calendars	leap
Ethiopic	Ethiopic	alms	Babylonian

Table 7: Calculated values for “calendar”

airplane			
<i>MIM</i>	<i>EMIM</i>	χ^2	<i>Dice</i>
USAFE	USAFE	MiG	MiG
Hiu	Hiu	maneuverability	plane
Alor	Alor	canopy	altitude
Plegovini	Plegovini	motherships	jets
bellow	bellow	Thunderstreak	maneuverability
RandalSchwartz	RandalSchwartz	underwing	pilots
Ufology	Ufology	84F	canopy
jib	jib	wrinkling	jet
Zhaoguo	Zhaoguo	Filmsite	Aviation
fashionably	fashionably	Maneuver	fuselage

Table 8: Calculated values for “airplane”

ocean			
<i>MIM</i>	<i>EMIM</i>	χ^2	<i>Dice</i>
Cheiro	Cheiro	Anstey	Antarctic
Tracysurf	Tracysurf	Bruticus	sail
Alvarolima	Alvarolima	DMeyering	floating
Dejima	Dejima	adverb	biodiversity
Sennet	Sennet	Paukrus	Fishing
iro	iro	tusk	ecosystems
Rockheights	Rockheights	bodyboarding	oceans
barque	barque	Orinoco	locked
bellow	bellow	plankton	temporarily
Ryanjunk	Ryanjunk	shack	seal

Table 9: Calculated values for “ocean”

bicycle			
<i>MIM</i>	<i>EMIM</i>	χ^2	<i>Dice</i>
Sergeants	Sergeants	racer	racer
Backhuys	Backhuys	cyclists	cyclists
Moetus	Moetus	Palmar��s	Discipline
Spudders	Spudders	Drunst	Palmar��s
Spilsby	Spilsby	Discipline	Drunst
Dockx	Dockx	Giro	cycling
MountainBikes	MountainBikes	U23	Giro
Klostergaard	Klostergaard	ProTeam	Friis
Lengerhane	Lengerhane	Friis	UCI
Khari	Khari	cycling	Rider

Table 10: Calculated values for “bicycle”

6 Appendix

```
1 #!/usr/bin/env python
2
3 import collections
4 import itertools
5 from data import words
6 from nltk.stem import *
7
8 class Result(object):
9     def __init__(self, a, b):
10         self.a = a
11         self.b = b
12         sa = set(words[a])
13         sb = set(words[b])
14         sab = sa.intersection(sb)
15         na = float(len(sa))
16         nb = float(len(sb))
17         nab = float(len(sab))
18         self.dice = nab / (na + nb)
19
20     def getdice(self):
21         return self.dice
22
23     def __repr__(self):
24         return '({},{}) Dice {}'.format(self.a, self.b, self.dice)
25
26 # skipping first 13,000 terms because they are all numeric
27 # so they won't meet language probability expectations
28 first1k = sorted(words.keys())[13000:14000]
29
30 # stem the first 1k words
31 stemmer = SnowballStemmer('english')
32 stems = {word: stemmer.stem(unicode(word, 'utf-8')) for word in first1k}
33
34 # count the stems to find duplicates
35 vals = collections.Counter(stems.values())
36
37 # reduce stem map to those that stemmed to the same stem
38 dupkeys = {key: val for key, val in stems.items() if vals[val] > 1}
39
40 # create new map that is the stem pointing to all terms that stemmed to it
41 dupset = {}
42 for pair in itertools.combinations(dupkeys.items(), 2):
43     k1 = pair[0][0]
44     k2 = pair[1][0]
45     v1 = pair[0][1]
46     v2 = pair[1][1]
47     if v1 == v2:
48         if not dupset.has_key(v1):
49             dupset[v1] = set()
50             dupset[v1].add(k1)
51             dupset[v1].add(k2)
52 print '%d duplicate stems' % len(dupset)
53
54 # calculate Dice's coefficient for each term with the same stem
55 results = {}
56 for stem, terms in dupset.items():
57     for pair in itertools.combinations(terms, 2):
58         t1 = pair[0]
59         t2 = pair[1]
60         if not results.has_key(stem):
61             results[stem] = set()
62             results[stem].add(Result(t1, t2))
```

Listing 2: stem.py

```

1 import cPickle
2
3 try:
4     print 'loading cached word map'
5     words = cPickle.load(open('words.p', 'rb'))
6 except IOError:
7     words = {line.split()[0]: line.split()[1:] for line in open('invidx.dat').readlines()}
8     cPickle.dump(words, open('words.p', 'wb'))
9     words = cPickle.load(open('words.p', 'rb'))
10 N = float(sum(len(docs) for docs in words.values()))

```

Listing 3: data.py

```

1 #!/usr/bin/env python
2
3 import re
4 import sys
5 import cPickle
6 from collections import Counter
7
8
9 def get_words():
10     try:
11         return cPickle.load(open('words.p', 'rb'))
12     except IOError:
13         wordmap = Counter(re.findall(r'\w+', open('big.txt').read().lower()))
14         cPickle.dump(wordmap, open('words.p', 'wb'))
15         return cPickle.load(open('words.p', 'rb'))
16
17 words = get_words()
18 N = sum(words.values())
19
20
21
22 def exists(wordset):
23     return set([word for word in wordset if word in words])
24
25
26 def prob(word):
27     return float(words[word]) / float(N)
28
29
30 def edit1(w):
31     letters = 'abcdefghijklmnopqrstuvwxyz'
32     deletes = [w[:i]+w[i+1:] for i in range(len(w))]
33     transposes = [w[:i]+w[i+1]+w[i]+w[i+2:] for i in range(len(w)-1)]
34     replaces = [w[:i]+l+w[i+1:] for i in range(len(w)) for l in letters]
35     inserts = [w[:i]+l+w[i:] for i in range(len(w)+1) for l in letters]
36     return set(deletes + transposes + replaces + inserts)
37
38
39 def edit2(word):
40     e2 = [edit1(w) for w in edit1(word)]
41     return [item for sublist in e2 for item in sublist]
42
43
44 def parse(word):
45     return exists([word]) or exists(edit1(word)) or exists(edit2(word)) or [word]
46
47
48 def correct(word):
49     return max(parse(word), key=prob)
50
51
52 if __name__ == '__main__':
53     print correct(sys.argv[1])

```

Listing 4: spelling.py

```

1 import cPickle
2 import math
3
4
5 class Result(object):
6     def __init__(self, a, b):
7         """calculate MIM, EMIM, Chi-square, and Dice's coefficient for words a and b.
8         mim = nab / (na * nb)
9         emim = nab * log [ N * nab / ( na * nb ) ]
10        x2 = ( nab - ( 1 / N ) * na * nb )^2 / ( na * nb )
11        dice = nab / ( na + nb )"""
12        self.a = a
13        self.b = b
14        sa = set(words[a])
15        sb = set(words[b])
16        sab = sa.intersection(sb)
17        na = float(len(sa))
18        nb = float(len(sb))
19        nab = float(len(sab))
20        self.mim = nab / (na * nb)
21        try:
22            self.emim = nab * math.log(N * nab / (na * nb))
23        except Exception as e:
24            self.emim = 0.0
25        self.x2 = (nab - (1/N) * na * nb)**2 / (na * nb)
26        self.dice = nab / (na + nb)
27
28    def getmim(self):
29        return self.mim
30
31    def getemim(self):
32        return self.emim
33
34    def getx2(self):
35        return self.x2
36
37    def getdice(self):
38        return self.dice
39
40    def __repr__(self):
41        return '{},{ }\n MIM {} \n EMIM {} \n X2 {} \n Dice {}'.format(
42            self.a, self.b, self.mim, self.emim, self.x2, self.dice)
43
44
45 def init():
46     global words
47     try:
48         print 'loading cached word map'
49         words = cPickle.load(open('words.p', 'rb'))
50     except IOError:
51         print 'cached word map not found, building now'
52         words = {line.split()[0]: line.split()[1:] for line in open('invidx.dat').readlines()
53                 }
54         cPickle.dump(words, open('words.p', 'wb'))
55         words = cPickle.load(open('words.p', 'rb'))
56     global N
57     N = float(sum(len(docs) for docs in words.values()))
58
59 def calc(choices):
60     print 'calculating...'
61     return {choice: [Result(choice, word) for word in words.keys() if choice != word] for
62             choice in choices}
63
64 def getthighest(results, choice, keyfunc):
65     return sorted(results[choice], key=keyfunc, reverse=True)[:10]
66
67
68 def printresults(results, choices):
69     print 'writing tables.tex'
70     with open('tables.tex', 'wb') as outfile:
71         for choice in choices:
72             mim = [res.b for res in getthighest(results, choice, Result.getmim)]
73             emim = [res.b for res in getthighest(results, choice, Result.getemim)]
74             x2 = [res.b for res in getthighest(results, choice, Result.getx2)]

```

```

75         dice = [res.b for res in getthighest(results, choice, Result.getdice)]
76         printtab(outfile, choice, mim, emim, x2, dice)
77
78
79 head = """\begin{table}[h!]
80 \centering
81 \begin{tabular}{l | c | c | c }
82 \hline
83 """
84
85 foot = '\hline\n\\end{tabular}\n\\caption{Calculated values for ‘%s\’}\n\\label{tab:
86         words}\n\\end{table}\n’
87
88 def printtab(outfile, choice, mim, emim, x2, dice):
89     outfile.write(head)
90     outfile.write('\multicolumn{4}{c}{’
91         + choice + ’}\n\\hline\n\\textit{MIM} & \\textit{EMIM} & \\textit{(\chi^2)} &
92         \\textit{Dice}\n\\hline\n’
93     for i in range(10):
94         outfile.write(row(i, mim, emim, x2, dice))
95     outfile.write(foot % choice)
96
97 def row(r, mim, emim, x2, dice):
98     return mim[r] + ’ & ’ + emim[r] + ’ & ’ + x2[r] + ’ & ’ + dice[r] + ’\n\\hline\n’
99
100 init()
101 choices = [
102     'running',
103     'calculation',
104     'color',
105     'horse',
106     'sky',
107     'railroad',
108     'calendar',
109     'airplane',
110     'ocean',
111     'bicycle']
112 results = calc(choices)
113 printresults(results, choices)

```

Listing 5: calc.py

7 References

- [1] Peter Norvig. How to Write a Spelling Corrector. Available at: <http://norvig.com/spell-correct.html>. Accessed: 2016/11/08.
- [2] The Python Programming Language. Available at: <https://www.python.org/>. Accessed: 2016/09/17.
- [3] Python.org. Python object serialization. Available at: <https://docs.python.org/2/library/pickle.html>. Accessed: 2016/11/06.