

Assignment 3

CS834-F16: Introduction to Information Retrieval

Fall 2016

Erika Siregar

CS Department - Old Dominion University

November 10, 2016

Question 6.1

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.

Answer

For indexing the collection, I modified the code that I used in assignment 2, sort the index alphabetically, and take the first 1000 words. The code for this is written in the file '1_2_index.py'. For creating the stem classes, I use Krovetz Stemmer because it produces a better stemming results than Porter Stemmer. Listing 1 shows the code to create the stem classes that utilize python library for Krovetz Stemmer [1].

```
1  #!/usr/bin/python
2  import json
3  from krovetzstemmer import Stemmer as KrovetzStemmer
4  import unicodedcsv as csv
5  from prettyprint import prettyprint
6
7
8  # Instantiate krovetz stemmer
9  krovetz = KrovetzStemmer()
10
11
12 # Read result of 1_index
13 with open('1_2_index.txt', 'rb') as f:
14     str_word_files_index = f.read()
15     word_files_index = json.loads(str_word_files_index)
16
17     stem_word_index = {}
18     for word, files in word_files_index.items():
19         # Stem word using krovetz
20         stemmed_word = krovetz.stem(word)
21
22         # Group by stemmed word
23         stem_word_index.setdefault(stemmed_word, [])
24         stem_word_index[stemmed_word].append(word)
25
26
27     filename = '3_stemmed_words.csv'
28     with open(filename, 'wb') as f:
29         print('Writing to file {}'.format(filename))
30
31         writer = csv.writer(f)
```

```

32     for stemmed_word, words in stem_word_index.items():
33         writer.writerow((stemmed_word, ', '.join(words)))
34
35     print('Done!')
```

Listing 1: Creating Stem Classes with Krovetz Stemmer

Table 1 shows the snippet of the stem classes created. The complete list of the stem classes is available in ‘3_stemmed_words.csv’ which is uploaded on github.

stem class	terms
academician	academicians, academician
adamant	adamantly, adamant
abundance	abundance
account	account, accounted, accounts, accounting
abdelkader	abdelkader
achter	achter
abednego	abednego
abortion	abortion, abortions
aboot	aboot
abrahamsson	abrahamsson
abdeali	abdeali
abandon	abandonment, abandon, abandoning

Table 1: A snippet of the stem classes

Next step is to create the stem clusters using Dice’s Coefficient [2] as the term association measure. Dice’s Coefficient works based on this formula:

$$2. \frac{n_{ab}}{n_a + n_b}$$

Using this formula, compute the Dice’s Coefficient for each pair of terms in every stem classes. With a threshold = 0.01, create a graph in which every pair of terms that has Dice’s Coefficient greater than 0.01 will be connected with an edge. Figure 1 shows the graph for the stem class ‘activate’ that can be grouped into two clusters: ‘*activate, activating, activator*’ and ‘*activates, activation*’. The graph for other stem classes are available on github in a folder named ‘*graph*’. From this graph, we only need to extract the connected components to form the clusters.

Listing 3 shows the code used to compute the Dice’s Coefficient, create the graphs, and extract the connected components of the graphs.

```

1
2  #!/usr/bin/python
3  import json
4  import nltk as nltk
5  from tabulate import tabulate
6  import unicodedcsv as csv
7  from pprint import pprint
8  import networkx as nx
9  import matplotlib.pyplot as plt
10
11  dice_coef_threshold = 0.01
12  stem_clusters = []
13
```

```

14 # Read result of 1_2_index.txt
15 with open('1_2_index.txt', 'rb') as f1:
16     word_files_index = json.loads(f1.read())
17
18 # Read result of 3_stemmed_words.csv
19 with open('3_stemmed_words.csv', 'rb') as f3:
20     for stemmed_word, words in csv.reader(f3):
21         words = words.split(',')
22
23         # create bigrams from words
24         bigrams = list(nltk.bigrams(words))
25         for word_a, word_b in bigrams:
26             # Lookup filename in word_files_index
27             files_a = word_files_index[word_a]
28             files_b = word_files_index[word_b]
29             files_a_sliced_b = list(set(files_b) & set(files_a))
30
31             dice_coef = float(2 * len(files_a_sliced_b)) / (len(files_a) + len(
32                 files_b))
33
34             if(dice_coef > dice_coef_threshold):
35                 stem_clusters.append((stemmed_word, word_a, word_b, dice_coef))
36
37 stem_clusters = sorted(stem_clusters, key=lambda x: x[3], reverse=True)
38 # print tabulate(stem_clusters, headers=['stemmed_word', 'word_a', 'word_b', '
39     dice_coef'])
40 filename = '4_dice_coefficient.csv'
41 with open(filename, 'wb') as f:
42     print('Writing to file {}'.format(filename))
43
44     writer = csv.writer(f)
45     for stemmed_word, word_a, word_b, dice_coef in stem_clusters:
46         writer.writerow((stemmed_word, word_a, word_b, dice_coef))
47
48
49 # Create graph
50 stemmed_word_data = {}
51 for stemmed_word, word_a, word_b, dice_coef in stem_clusters:
52     stemmed_word_data.setdefault(stemmed_word, [])
53     stemmed_word_data[stemmed_word].append((word_a, word_b, dice_coef))
54
55 stemmed_word_clusters = {}
56 for stemmed_word, data in stemmed_word_data.items():
57     G=nx.MultiGraph()
58
59     labels = {}
60     for word_a, word_b, dice_coef in data:
61         G.add_edge(word_a, word_b, weight=dice_coef, label=dice_coef)
62         labels[(word_a, word_b)] = dice_coef
63
64 # export connected components into list
65 stemmed_word_clusters[stemmed_word] = list(nx.connected_components(G))
66
67 nx.draw(G, with_labels=True)
68 nx.draw_networkx_edge_labels(G, pos=nx.spring_layout(G), edge_labels=labels)
69
70 filename = '4_graph_{}.png'.format(stemmed_word)

```

```

71     print('Saving graph {}'.format(filename))
72     plt.savefig(filename, format='PNG')
73     plt.clf()
74
75     print('Draw graphics done!')
76
77     print('Print stem clusters...')
78
79     for stemmed_word, connected_nodes in stemmed_word_clusters.items():
80         for connected_node in connected_nodes:
81             print(u'{}\t: {}'.format(stemmed_word, ', '.join(connected_node)))
82
83     print('Print stem clusters done')

```

Listing 2: Creating Cluster using Dice's Coefficient

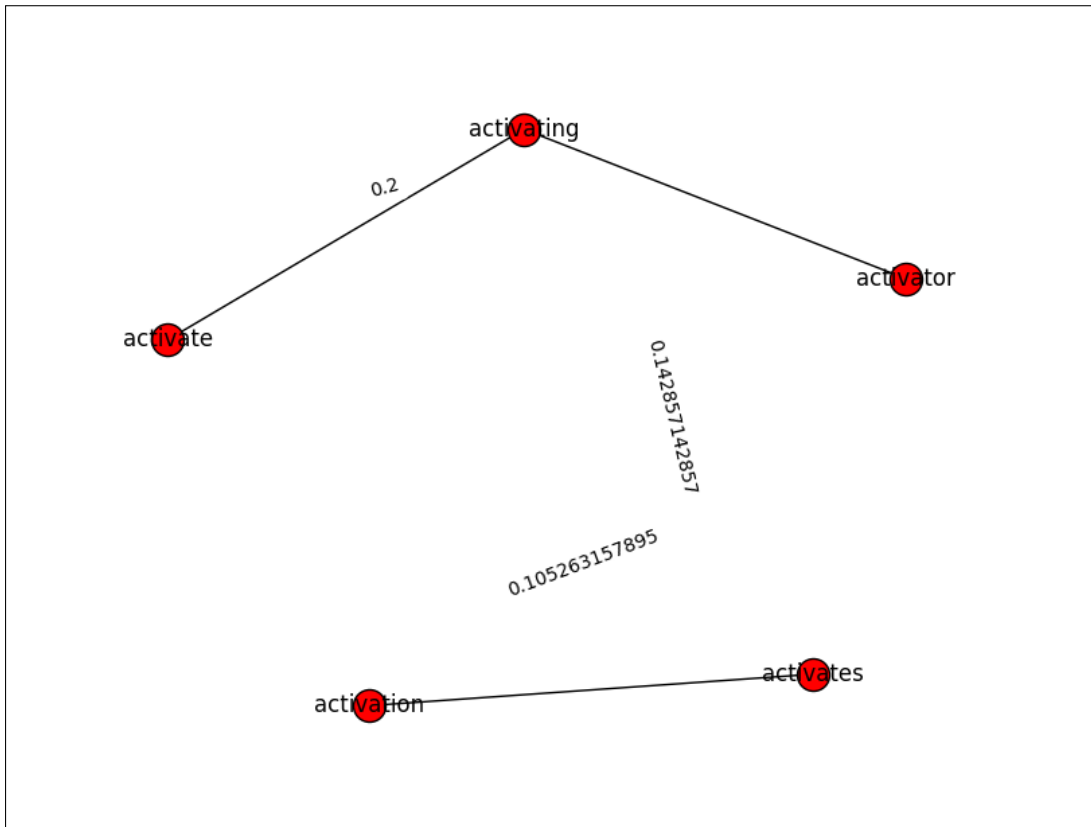


Figure 1: Graph of the connected component for the stem class 'activate'

Table 2 shows the Dice's Coefficient for some pair of terms in the Wikipedia collection. The complete list of the term pairs with their Dice's Coefficient value is available on Github in a file named '4_dice_coefficient.csv'.

No	Class	Term 1	Term 2	Dice's Coefficient
1	abomination	abominations	abomination	0.6666666667
2	abut	abuts	abutting	0.6666666667
3	abjad	abjads	abjad	0.6666666667
4	academician	academicians	academician	0.5714285714
5	aberration	aberrations	aberration	0.5
6	adapter	adapters	adapter	0.5
7	actor	actors	actor	0.4984025559
8	abridge	abridged	abridges	0.4
9	absolve	absolve	absolved	0.4
10	acoustic	acoustically	acoustical	0.4

Table 2: Dice's Coefficient for some pair of terms in small Wikipedia collection

Listing 3 shows the clusters resulted for the small Wikipedia Collection.

```

1 accessible : accessible , accessibility
2 activate : activator , activate , activating
3 activate : activation , activates
4 accelerator : accelerators , accelerator
5 abridge : abridges , abridged
6 accurate : accurate , accurately
7 address : addressing , addresses , addressed , address
8 abomination : abomination , abominations
9 accept : accepting , accepted , accept
10 abbreviation : abbreviation , abbreviations
11 acclaim : acclaim , acclaimed
12 adaptation : adaptations , adaptation
13 abrogate : abrogation , abrogated
14 accrete : accreting , accrete , accreted
15 acid : acidic , acids , acid
16 accommodate : accommodate , accommodated
17 absolute : absolute , absolute
18 acceleration : acceleration , accelerations
19 additional : additional , additionally
20 acknowledge : acknowledges , acknowledged
21 addition : addition , additions
22 accent : accent , accented
23 actor : actors , actor
24 access : access , accessed , accessing , accessor
25 acyltransferase : acyltransferase , acyltransferases
26 add : adding , add , added , adds
27 activist : activist , activists
28 adapt : adaption , adapt
29 adapt : adaptive , adapted
30 acre : acres , acre
31 achieve : achieves , achieve
32 achieve : achieving , achieved
33 abstraction : abstraction , abstractions
34 accompany : accompany , accompanying
35 activity : activities , activity
36 accidental : accidentally , accidental
37 aberration : aberrations , aberration
38 acronym : acronym , acronyms
39 academy : academy , academies
40 acquire : acquire , acquiring

```

```

41 academician : academician, academicians
42 abut : abutting, abuts
43 abuse : abused, abuse
44 accompaniment : accompaniment, accompaniments
45 actress : actresses, actress
46 accuse : accusing, accuses
47 acute : acute, acutely
48 accumulate : accumulate, accumulated
49 abugida : abugidas, abugida
50 abduct : abducted, abductors
51 achievement : achievements, achievement
52 accredit : accrediting, accredited, accreditation
53 accusation : accusation, accusations
54 account : accounting, accounted, accounts
55 accident : accident, accidents
56 actual : actual, actualized
57 adapter : adapter, adapters
58 accomplishment : accomplishment, accomplishments
59 absolve : absolved, absolve
60 abjad : abjads, abjad
61 academic : academic, academically
62 abrupt : abrupt, abruptly
63 abolition : abolitionism, abolition
64 act : acted, act
65 action : action, actions
66 acoustic : acoustical, acoustic, acoustically
67 abbey : abbeys, abbey
68 acquisition : acquisitions, acquisition
69 abbot : abbots, abbot

```

Listing 3: Cluster for the small wikipedia collection

Question 6.2

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).

Answer:

Spelling corrector based on the noisy channel model works using this approach [3]:

1. Let's say x is a misspelled word and $w = w_1, w_2, w_3, \dots, w_n$ is an array of possible corrected words.
2. Our task is to compute the conditional probability and take a word w_i that has the maximum value for $P(x|w_i) \cdot P(w_i)$.

It looks complicated. But, fortunately, Peter Norvig [4] has provided a nice python code for spelling corrector, which is worked based on the noisy channel model. To compute the probability, Norvig uses word frequencies taken from a predefined dataset 'big.txt'. For this assignment, I modified Norvig's code as can be seen in listing 4.

```

1
2 #!/usr/bin/python
3
4 import re
5 from collections import Counter
6
7 import sys
8
9
10 class Spell:
11     def __init__(self, train_file):
12         self.document = open(train_file).read()
13         self.to_words()
14         self.count_words()
15
16     def to_words(self):
17         self.words = re.findall(r'\w+', self.document.lower())
18
19     def count_words(self):
20         self.word_count = Counter(self.words)
21
22     def probability(self, word):
23         "Probability of 'word'."
24         word = word.lower()
25         N = sum(self.word_count.values())
26         return self.word_count[word] / N
27
28     def edits1(self, word):
29         "All edits that are one edit away from 'word'."
30         word = word.lower()
31         letters = 'abcdefghijklmnopqrstuvwxyz'
32         splits = [(word[:i], word[i:]) for i in range(len(word) + 1)]
33         deletes = [L + R[1:] for L, R in splits if R]
34         transposes = [L + R[1] + R[0] + R[2:] for L, R in splits if len(R) > 1]
35         replaces = [L + c + R[1:] for L, R in splits if R for c in letters]
36         inserts = [L + c + R for L, R in splits for c in letters]
37         return set(deletes + transposes + replaces + inserts)
38
39     def edits2(self, word):
40         "All edits that are two edits away from 'word'."
41         word = word.lower()
42         return {e2 for e1 in self.edits1(word) for e2 in self.edits1(e1)}
43
44     def known(self, words):
45         "The subset of 'words' that appear in the dictionary of WORDS."
46         return set(w for w in words if w in self.word_count)
47
48     def candidates(self, word):
49         "Generate possible spelling corrections for word."
50         word = word.lower()
51         return (self.known([word]) or self.known(self.edits1(word)) or self.known(
52 self.edits2(word)) or [word])
53
54     def correction(self, word):
55         "Most probable spelling correction for word."
56         word = word.lower()
57         candidates = self.candidates(word)
58         return max(candidates, key=self.probability)

```



```

59
60 if __name__ == '__main__':
61     spell = Spell('big.txt')
62
63     if len(sys.argv) != 2:
64         print('python spell.py <word>')
65         exit()
66
67     word = sys.argv[1]
68     spelled_word = spell.correction(word)
69
70     print('{} --> {}'.format(word, spelled_word))

```

Listing 4: Spelling corrector

Figure 2 shows the example of spelling correction using some words taken from the textbook [5] and some words taken from Norvig’s website [4].

```

erikaris@inspiron-7368:/media/erikaris/DATA/ODU/Semester 3/intro_to_info_retrieval/assignments/a3/6_2$ python spell.py speling
speling --> spelling
erikaris@inspiron-7368:/media/erikaris/DATA/ODU/Semester 3/intro_to_info_retrieval/assignments/a3/6_2$ python spell.py korrekstud
korrekstud --> corrected
erikaris@inspiron-7368:/media/erikaris/DATA/ODU/Semester 3/intro_to_info_retrieval/assignments/a3/6_2$ python spell.py extensions
extensions --> extensions
erikaris@inspiron-7368:/media/erikaris/DATA/ODU/Semester 3/intro_to_info_retrieval/assignments/a3/6_2$ python spell.py inconvient
inconvient --> inconvenient
erikaris@inspiron-7368:/media/erikaris/DATA/ODU/Semester 3/intro_to_info_retrieval/assignments/a3/6_2$ python spell.py poiner
poiner --> pointer
erikaris@inspiron-7368:/media/erikaris/DATA/ODU/Semester 3/intro_to_info_retrieval/assignments/a3/6_2$ python spell.py peotryy
peotryy --> poetry
erikaris@inspiron-7368:/media/erikaris/DATA/ODU/Semester 3/intro_to_info_retrieval/assignments/a3/6_2$ python spell.py brimingham
brimingham --> birmingham

```

Figure 2: Spelling corrector using noisy channel model with edit distance of 1 and 2

Question 6.5

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.

Answer

For this question, I downloaded the source code for Galago version 3.10 [6] from <https://sourceforge.net/p/lemur/galago/ci/release-3.10/tree/>. It is a part of The Lemur Project [7]. The code for the snippet generator can be found in a file named ‘SnippetGenerator.java’ under the directory ‘galago-3.10/core/src/main/java/org/lemurproject/galago/core/index/corpus’. This is how the code works:

1. Given ‘documentText’ and a set of ‘queryTerms’, tokenize the ‘documentText’ into terms and its position. Listing 5 shows the code for this step.
2. Stem each term using a defined stemmer. By default, Galago uses Krovetz stemmer.
3. Iterate each stemmed term in documentText and find matches with each term in queryTerm.

4. For each matched term, make snippet region containing match term (original term) and maximum 5 terms before and 4 terms after original term in document. So, the snippet region will contain maximum 10 terms including the matched term.
5. Check the snippet regions and resolve if there are overlapped regions.
6. Remove snippet regions that overflow the maxSize. In the code, the maxSize is set to 40.
7. Combine all snippet regions in each document into a single snippet splitted by ‘...’ (three dots). Make the matched terms displayed in a bold format.

```

1 public String getSnippet(String documentText, Set<String> queryTerms) throws
   IOException {
2     ArrayList<IntSpan> positions = new ArrayList<IntSpan>();
3     Document document = parseAsDocument(documentText, positions);
4     return generateSnippet(document, positions, queryTerms);
5 }

```

Listing 5: Tokenizing in Galago’s snippet generator

Figure ?? shows the example of running the Galago’s snippet generator on one of documents in the small Wikipedia collections ‘xxx.html’.

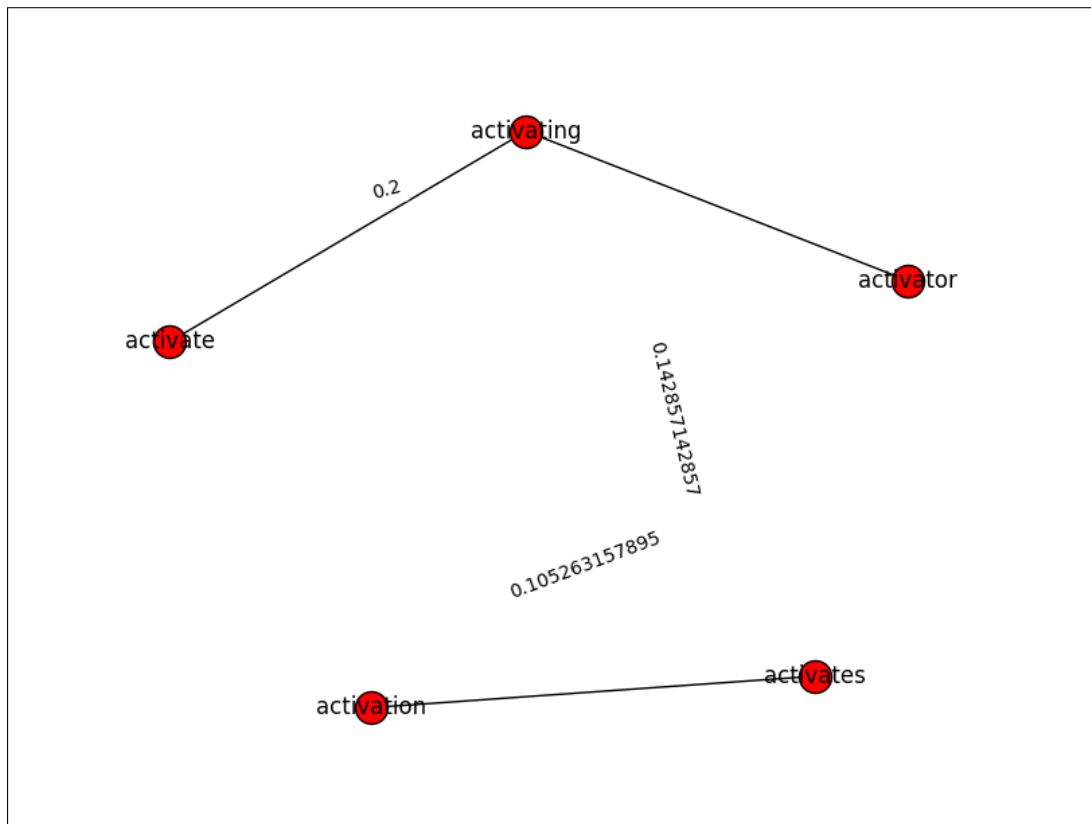


Figure 3: xxx

I also run the snippet generator on a document that has small content. For this assignment, I modified the document ‘xxy.html’ so that it only contains 2 sentences. Figure 4 shows the output of Galago’s snippet generator for a document with small content.

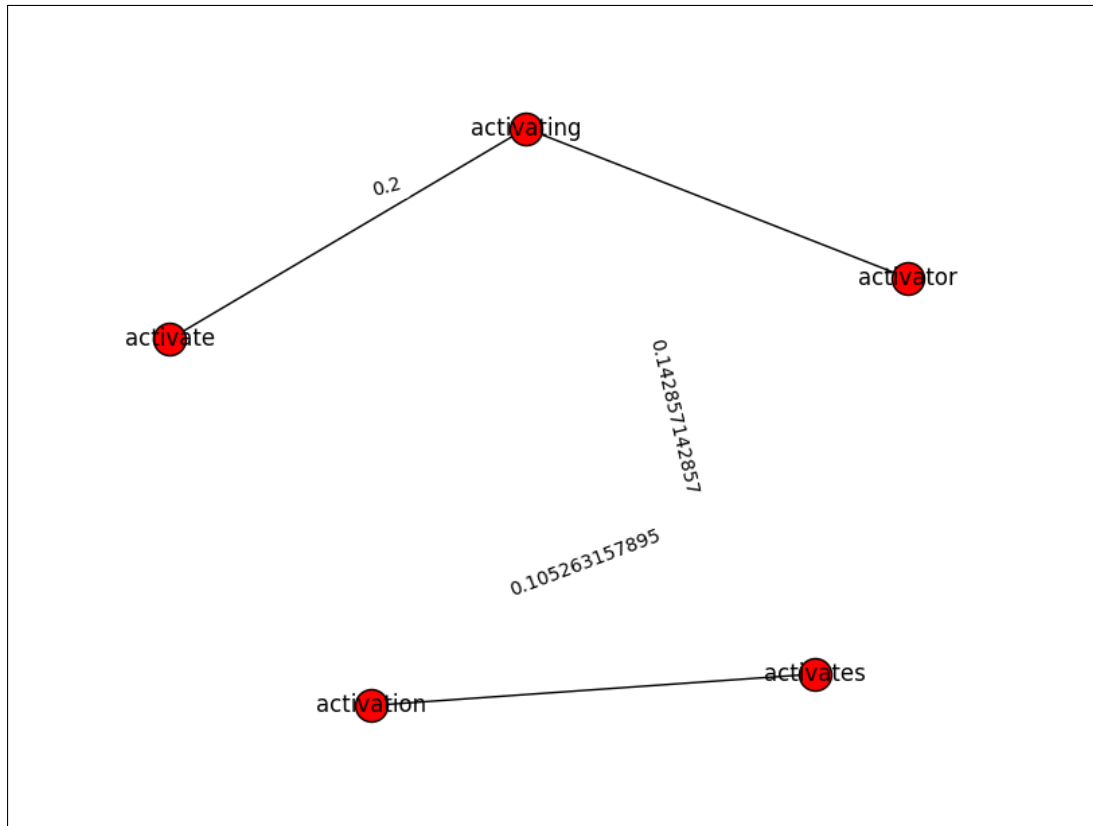


Figure 4: xxy

From figure 4, I conclude that ... for document with little content, the algorithm I also find out that the algorithm does not handle the phrase well. For example, if I type a query ‘United States’, it will return a snippet as can be seen in figure 5. The algorithm split the phrase and find the best match based on term-by-term, not based on the complete phrase (two terms at a time).

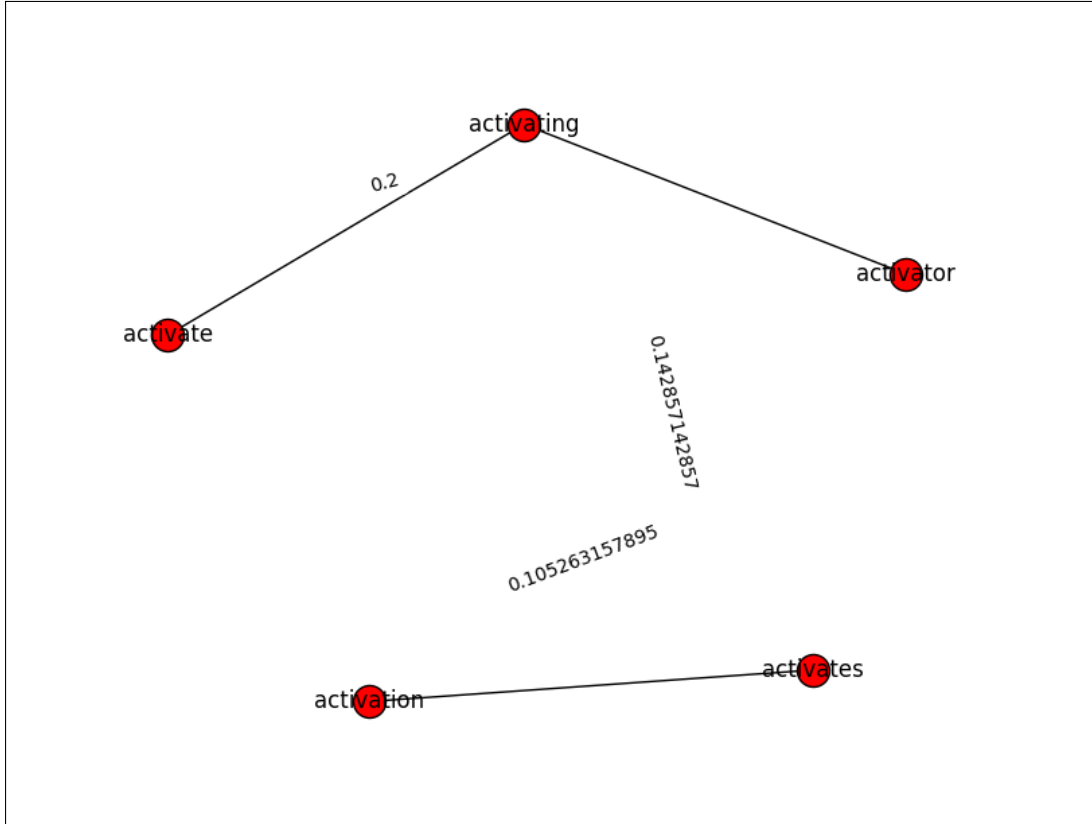


Figure 5: xxz

One thing that we can use to improve the algorithm is to adopt the algorithm explained by Turpin [8]. In his algorithm, Turpin [8] takes ‘the longest contiguous run of query terms’ as a variable to rank the sentences. Using this approach, I think we can avoid getting the snippet contains separated query terms. Figure 6

IN	A document broken into one sentence per line, and a sequence of query terms.
1	For each line of the text, $\mathcal{L} = [w_1, w_2, \dots, w_m]$
2	Let h be 1 if \mathcal{L} is a heading, 0 otherwise.
3	Let ℓ be 2 if \mathcal{L} is the first line of a document, 1 if it is the second line, 0 otherwise.
4	Let c be the number of w_i that are query terms, counting repetitions.
5	Let d be the number of distinct query terms that match some w_i .
6	Identify the longest contiguous run of query terms in \mathcal{L} , say $w_j \dots w_{j+k}$.
7	Use a weighted combination of c, d, k, h and ℓ to derive a score s .
8	Insert \mathcal{L} into a max-heap using s as the key.
OUT	Remove the number of sentences required from the heap to form the summary.

Figure 6: Turpin's algorithm for ranking the sentences. Adapted from [8]

Question MLN1

MLN1: using the small wikipedia example, choose 10 words and create stem classes as per the algorithm on pp. 191-192.

Answer

D

Question MLN2

MLN2: 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)

Answer

e

References

- [1] Ruey-Cheng Chen. KrovetzStemmer 0.4. <https://pypi.python.org/pypi/KrovetzStemmer/0.4>, 2016. [Online; accessed 12-October-2016].
- [2] Wikipedia. Algorithm Implementation/Strings/Dice's coefficient. https://en.wikibooks.org/wiki/Algorithm_Implementation/Strings/Dice%27s_coefficient/, 2016. [Online; accessed 8-November-2016].
- [3] Dan Jurafsky. Spelling Correction and the Noisy Channel . <https://web.stanford.edu/class/cs124/lec/spelling.pdf>, 2012. [Online; accessed 8-November-2016].
- [4] Peter Norvig. How to Write a Spelling Corrector . <http://norvig.com/spell-correct.html>, 2016. [Online; accessed 9-November-2016].
- [5] Bruce Croft, Donald Metzler, and Trevor Strohman. *Search Engines: Information Retrieval in Practice*. Addison-Wesley Publishing Company, USA, 1st edition, 2009.
- [6] The Lemur Project. The Lemur Project - Galago 3.10 . <https://sourceforge.net/p/lemur/galago/ci/release-3.10/tree/>, 2016. [Online; accessed 5-November-2016].
- [7] The Lemur Project. The Lemur Project . <http://www.lemurproject.org/>, 2016. [Online; accessed 5-November-2016].
- [8] Andrew Turpin, Yohannes Tsegay, David Hawking, and Hugh E. Williams. Fast generation of result snippets in web search. In *Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '07, pages 127–134, New York, NY, USA, 2007. ACM.