

Assignment 3

Fall 2016

CS834 Introduction to Information Retrieval

Dr. Michael Nelson

Mathew Chaney

November 6, 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 | Appendix | 7 |
| 6 | References | 9 |

List of Figures

Listings

| | | |
|---|--------------------------------------|---|
| 1 | spelling.py example output | 4 |
| 2 | spelling.py | 7 |
| 3 | stem.py | 8 |

List of Tables

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 3, 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

The `spelling.py` script, found in Listing 2, was created using the Python programming language [1]. Downloaded `big.txt` to calculate an example language model. These words were counted and stored in a map that was compressed on disk using the pickle python library [2].

The process of determining a spelling correction is as follows:

1. Count all the words contained in the example text file. This count is used to determine the language model $P(W)$ calculation.
2. Take the input word and determine all existing (correctly spelled) words with edit distance one and two.
3. With the assumption that shorter edit distances equate to a higher probability of being the correct spelling, select from the remaining set of (valid) words the one with the shortest edit distance and highest value for $P(W)$.

$P(W)$ 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.

2.3 Results

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

```
1 [mchaney@mchaney-1 spelling]$ ./spelling.py words
2 words
3 [mchaney@mchaney-1 spelling]$ ./spelling.py flewurz
4 flower
5 [mchaney@mchaney-1 spelling]$ ./spelling.py spilling
6 selling
7 [mchaney@mchaney-1 spelling]$ ./spelling.py swelling
8 swelling
9 [mchaney@mchaney-1 spelling]$ ./spelling.py aacck
10 back
```

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 Appendix

```
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 2: spelling.py

```
1 #!/usr/bin/env python
```

Listing 3: stem.py

6 References

- [1] The Python Programming Language. Available at: <https://www.python.org/>. Accessed: 2016/09/17.
- [2] Python.org. Python object serialization. Available at: <https://docs.python.org/2/library/pickle.html>. Accessed: 2016/11/06.