# 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 alpha-
betically, 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].

```python
#!/usr/bin/python
import json
from krovetzstemmer import Stemmer as KrovetzStemmer
import unicodecsv as csv
from prettyprint import prettyprint


# Instantiate krovetz stemmer
krovetz = KrovetzStemmer()


# Read result of 1_index
with open('1_2_index.txt', 'rb') as f:
    str_word_files_index = f.read()
    word_files_index = json.loads(str_word_files_index)

    stem_word_index = {}
    for word, files in word_files_index.items():
        # Stem word using krovetz
        stemmed_word = krovetz.stem(word)

        # Group by stemmed word
        stem_word_index.setdefault(stemmed_word, [])
        stem_word_index[stemmed_word].append(word)


    filename = '3_stemmed_words.csv'
    with open(filename, 'wb') as f:
        print('Writing to file {}'.format(filename))

        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:

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

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 unicodecsv as csv
7  from prettyprint import prettyprint
8  import networkx as nx
9  import matplotlib.pyplot as plt
10
11 dice_coef_threshold = 0.01
12 stem_clusters = []
13
```

```python
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(
    files_b))
32
33                  if(dice_coef > dice_coef_threshold):
34                      stem_clusters.append((stemmed_word, word_a, word_b, dice_coef))
35
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', '
    dice_coef'])
39
40  filename = '4_dice_coeficient.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')
```

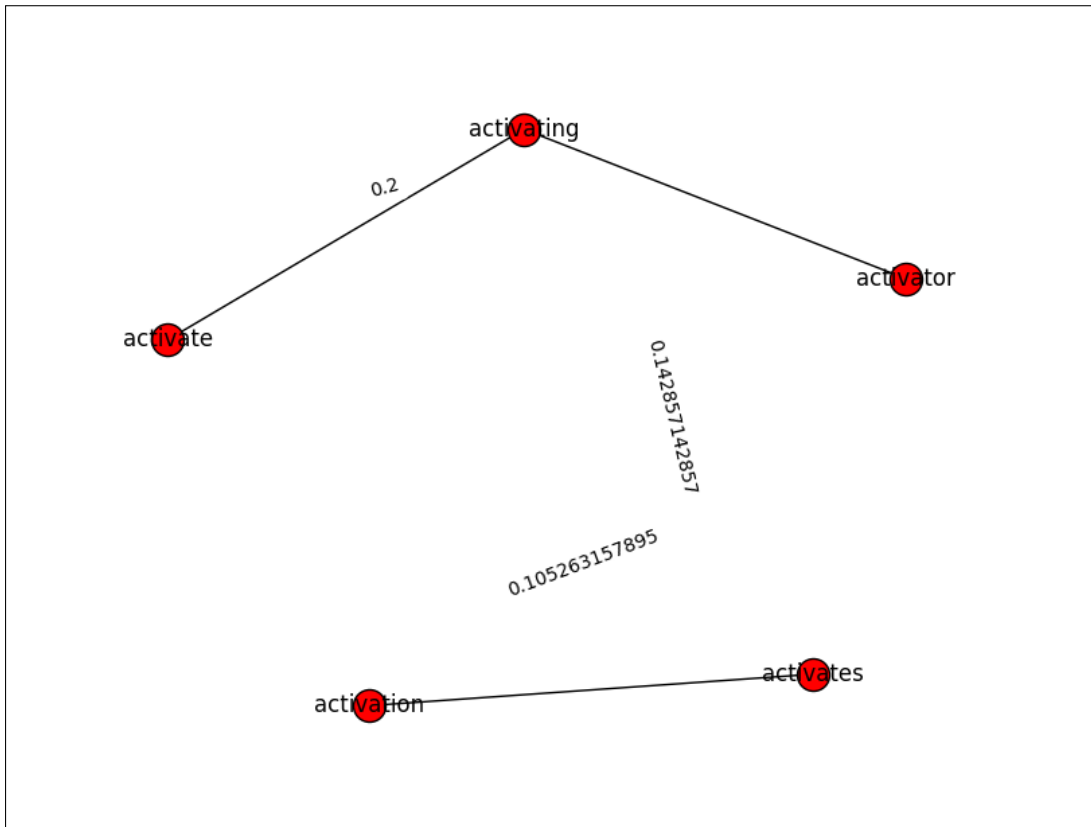Listing 2: Creating Cluster using Dice's Coefficient



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_coeficient.csv'.

| No | Class | Term 1 | Term 2 | Dice's Coefficient |
|---|---|---|---|---|
| 1 | abomination | abominationes | 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, abominationes
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   : absoluter, 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 mispelled word and $w = w_1, w_2, w_3, ..., w_n$ 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 $P(x|w_i) . P(w_i)$.

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.

```python
#!/ usr / bin / python

import re
from collections import Counter

import sys


class Spell :
    def __init__ ( self , train_file ) :
        self . document = open ( train_file ) . read ()
        self . to_words ()
        self . count_words ()

    def to_words ( self ) :
        self . words = re . findall ( r '\w+', self . document . lower () )

    def count_words ( self ) :
        self . word_count = Counter ( self . words )

    def probability ( self , word ) :
        " Probability of 'word'. "
        word = word . lower ()
        N = sum ( self . word_count . values () )
        return self . word_count [ word ] / N

    def edits1 ( self , word ) :
        " All edits that are one edit away from 'word'. "
        word = word . lower ()
        letters = 'abcdefghijklmnopqrstuvwxyz'
        splits = [ ( word [: i ] , word [ i :] ) for i in range ( len ( word ) + 1 ) ]
        deletes = [ L + R [1:] for L , R in splits if R ]
        transposes = [ L + R [1] + R [0] + R [2:] for L , R in splits if len (R) > 1 ]
        replaces = [ L + c + R [1:] for L , R in splits if R for c in letters ]
        inserts = [ L + c + R for L , R in splits for c in letters ]
        return set ( deletes + transposes + replaces + inserts )

    def edits2 ( self , word ) :
        " All edits that are two edits away from 'word'. "
        word = word . lower ()
        return ( e2 for e1 in self . edits1 ( word ) for e2 in self . edits1 ( e1 ) )

    def known ( self , words ) :
        " The subset of 'words' that appear in the dictionary of WORDS. "
        return set (w for w in words if w in self . word_count )

    def candidates ( self , word ) :
        " Generate possible spelling corrections for word. "
        word = word . lower ()
        return ( self . known ([ word ]) or self . known ( self . edits1 ( word )) or self . known (
    self . edits2 ( word )) or [ word ])

    def correction ( self , word ) :
        " Most probable spelling correction for word. "
        word = word . lower ()
        candidates = self . candidates ( word )
        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].



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 1 shows the code for this step.

```
1
2  public String getSnippet(String documentText, Set<String> queryTerms) throws
       IOException {
```

```
3        ArrayList<IntSpan> positions = new ArrayList<IntSpan>();
4        Document document = parseAsDocument(documentText, positions);
5        return generateSnippet(document, positions, queryTerms);
6    }
```

2. Stem each term using a defined stemmer. By default, Galago uses Krovetz stemmer.

```
1
2    private Document parseAsDocument(String text, ArrayList<IntSpan> positions)
         throws IOException {
3        Document document = new Document();
4        document.text = text;
5
6        // Tokenize the document
7        TagTokenizer tokenizer = new TagTokenizer();
8        tokenizer.process(document);
9
10       if (positions != null) {
11           positions.addAll(tokenizer.getTokenPositions());
12       }
13       if (stemming) {
14           document = stemmer.stem(document);
15       }
16
17       return document;
18   }
```

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.

```
1
2    Private ArrayList<SnippetRegion> findMatches(final Document document, final Set
         <String> queryTerms) {
3        // Make a snippet region object for each term occurrence in the document,
4        // while also counting matches
5        ArrayList<SnippetRegion> regions = new ArrayList<SnippetRegion>();
6
7        for (int i = 0; i < document.terms.size(); i++) {
8            String term = document.terms.get(i);
9            if (queryTerms.contains(term)) {
10               regions.add(new SnippetRegion(term, i, width, document.terms.size()));
11           }
12       }
13       return regions;
14   }
```

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.

```
1
2    public ArrayList<SnippetRegion> combineRegions(final ArrayList<SnippetRegion>
         regions) {
3        ArrayList<SnippetRegion> finalRegions = new ArrayList<SnippetRegion>();
4        SnippetRegion last = null;
5        int snippetSize = 0;
```

9

```
6        int maxSize = 40;
7
8        for (SnippetRegion current : regions) {
9          if (last == null) {
10           last = current;
11         } else if (last.overlap(current)) {
12           SnippetRegion bigger = last.merge(current);
13
14           if (bigger.size() + snippetSize > maxSize) {
15             finalRegions.add(last);
16             last = null;
17           } else {
18             last = bigger;
19           }
20         } else if (last.size() + snippetSize > maxSize) {
21           break;
22         } else {
23           finalRegions.add(last);
24           snippetSize += last.size();
25           last = current;
26         }
27       }
28
29       if (last != null && snippetSize + last.size() < maxSize) {
30         finalRegions.add(last);
31       }
32
33       return finalRegions;
34     }
```

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 buildHtmlString(Snippet best, Document document, ArrayList<
    IntSpan> positions) {
2     StringBuilder builder = new StringBuilder();
3
4     for (SnippetRegion region : best.regions) {
5       if (region.start != 0) {
6         builder.append("...");
7       }
8       int startChar = positions.get(region.start).start;
9       int endChar = positions.get(region.end - 1).end;
10      int start = 0;
11
12      // section string
13      String section = document.text.substring(startChar, endChar);
14
15      for (Match m : region.matches) {
16        int startMatchChar = positions.get(m.start).start - startChar;
17        int endMatchChar = positions.get(m.end - 1).end - startChar;
18
19        String intermediate = stripTags(section.substring(start, startMatchChar)
    );
20        builder.append(intermediate);
21        builder.append("<strong>");
22        builder.append(stripTags(section.substring(startMatchChar, endMatchChar)
    ));
23        builder.append("</strong>");
```

```
24            start = endMatchChar;
25          }
26
27          if (start >= 0) {
28            builder.append(stripTags(section.substring(start)));
29          }
30
31          // terminate matches once we reached a max length.
32          int maxSnippetSize = 500;
33          if (builder.length() > maxSnippetSize) {
34            break;
35          }
36        }
37
38        if (best.regions.size() > 1 && best.regions.get(best.regions.size() - 1).end
          != document.terms.
39                size()) {
40          builder.append("...");
41        }
42        return builder.toString();
43      }
```

Figure 3 shows the example of snippets generated for the query terms 'computer science'.



Figure 3: Snippet generated by Galago for query terms 'computer science'

Based on my analysis, I think Galago's snippet algorithm will not work well for pages with little text content. It is because the algorithm works by matching the query terms with the page content. If there is only little content on the page, then the probability to generate a good snippet will also decrease. One thing that we can use to improve the algorithm is to adopt the algorithm explained by Turpin [8]. Turpin's algorithm does not only depends on the 'matched terms'. It also uses the combination of several weights to improve the rank of the sentences such as the longest contiguous

run of query terms and the position of the sentences. Figure 4 illustrates the algorithm used by Turpin [8] to generate snippets.



| | |
|---|---|
| **IN** | A document broken into one sentence per line, and a sequence of query terms. |
| **1** | For each line of the text, $\mathcal{L} = [w_1, w_2, \ldots, w_m]$ |
| **2** | Let $h$ be 1 if $\mathcal{L}$ is a heading, 0 otherwise. |
| **3** | Let $\ell$ 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 $w_j \ldots w_{j+k}$. |
| **7** | Use a weighted combination of $c$, $d$, $k$, $h$ and $\ell$ 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 4: 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

The steps to solve this problem is similar to what I have done in question 6.1. There are some differences:

1. There are 4 metrics that are use to measure the association between the terms, which are Dice's Coefficient, EMIM, MIM, and Chi-square.

2. Window size = 100 words

3. Ten chosen words, which are *daughter, army, provide, type, feature, free, standard, street, fin, and unit.*

Listing 5 shows the python code for this problem.

```
#!/usr/bin/python
import io
import math
import os
```

```python
6  from krovetzstemmer import Stemmer as KrovetzStemmer
7
8  import html2text
9  import matplotlib.pyplot as plt
10 import networkx as nx
11 import nltk
12
13
14 # Create graph function
15 def create_graph(stem_data, metric_name):
16     stemmed_word_data = {}
17     for stemmed_word, word_a, word_b, coef in stem_data:
18         stemmed_word_data.setdefault(stemmed_word, [])
19         stemmed_word_data[stemmed_word].append((word_a, word_b, coef))
20
21     stemmed_word_clusters = {}
22     for stemmed_word, data in stemmed_word_data.items():
23         G = nx.MultiGraph()
24
25         labels = {}
26         for word_a, word_b, coef in data:
27             G.add_edge(word_a, word_b, weight=coef, label=coef)
28             labels[(word_a, word_b)] = coef
29
30         # export connected components into list
31         stemmed_word_clusters[stemmed_word] = list(nx.connected_components(G))
32
33         nx.draw(G, with_labels=True)
34         nx.draw_networkx_edge_labels(G, pos=nx.spring_layout(G), edge_labels=labels)
35
36         filename = '{}_graph_{}.png'.format(metric_name, stemmed_word)
37         print('Saving graph {}'.format(filename))
38         plt.savefig(filename, format='PNG')
39         plt.clf()
40
41     return stemmed_word_clusters
42
43
44 # Print stem cluster function
45 def print_stem_cluster(stemmed_word_clusters, metric_name):
46     print('\n')
47     print('Clusters using metric {}'.format(metric_name))
48     print('=================================')
49
50     for stemmed_word, connected_nodes in stemmed_word_clusters.items():
51         for connected_node in connected_nodes:
52             print(u'{}\t: {}'.format(stemmed_word, ', '.join(connected_node)))
53
54
55 # List all files
56 html_files = []
57 for root, dirs, files in os.walk(os.path.abspath('/media/erikaris/DATA/ODU/Semester
       3/intro_to_info_retrieval/assignments/a2/code_report/articles/z')):
58     for file in files:
59         if file.endswith('.html'):
60             filepath = os.path.join(root, file)
61             html_files.append(filepath)
```

```python
62
63
64 file_words_index = {}
65 all_words = set()
66 window_text_number = 100
67
68 # Index all words and files
       ================================================================
69 for idx, file in enumerate(html_files):
70     print('{} of {}. Processing file {}'.format(idx+1, len(html_files), file))
71     print('=' * 30)
72
73     # get text only from each file -> remove all tags
74     h = html2text.HTML2Text()
75     h.ignore_links = True
76     text = h.handle(u' '.join([line.strip() for line in io.open(file, "r", encoding="
       utf-8").readlines()]))
77
78     # get all words from text by splitting by whitespace
79     words = [word.lower() for word in text.split() if word.isalpha()]
80
81     # make text window containing 50 - 100 words
82     for w in range(1, int(len(words)/window_text_number)):
83         start = (w-1)*window_text_number
84         end = w*window_text_number
85         if end > len(words): end = len(words)
86
87         window_words = words[start:end]
88
89         all_words |= set(window_words)
90         file_words_index[os.path.basename(file) + '_part_' + str(w)] = window_words
91
92
93 # # Remove stopword to purify result
       ================================================================
94 # nltk.download('stopwords')
95 # all_words = [word for word in all_words if word not in stopwords.words('english')]
96 #
97 #
98 # # Select random 10 words
       ================================================================
99 # all_words = random.sample(all_words, 10)
100
101
102 # Invert words and files index
       ================================================================
103 word_files_index = {}
104 for idx, word in enumerate(all_words):
105     print('{} of {}. Processing word {}'.format(idx + 1, len(all_words), word.encode(
       'utf-8')))
106     print('=' * 30)
107
108     files = []
109     for file, words in file_words_index.items():
110         if word in words:
111             files.append(file)
112     word_files_index[word] = sorted(set(files))
113
```

```python
114
115  # Stem words
       =================================================================================================

116  krovetz = KrovetzStemmer()
117  stem_word_index = {}
118  for word, files in word_files_index.items():
119      # Stem word using krovetz
120      stemmed_word = krovetz.stem(word)
121
122      # Group by stemmed word
123      stem_word_index.setdefault(stemmed_word, [])
124      stem_word_index[stemmed_word].append(word)
125
126
127  # Calculate coefficient
       =================================================================================================

128  coef_threshold = 0.0
129
130  dice_stemmed_word_data = []
131  mim_stemmed_word_data = []
132  emim_stemmed_word_data = []
133  chi_sqr_stemmed_word_data = []
134
135  counter = 0
136  for stemmed_word, words in stem_word_index.items():
137      # create bigrams from words
138      bigrams = list(nltk.bigrams(words))
139      for word_a, word_b in bigrams:
140          # Lookup filename in word_files_index
141          files_a = word_files_index[word_a]
142          files_b = word_files_index[word_b]
143          files_a_sliced_b = list(set(files_b) & set(files_a))
144
145          # Using dice coef
146          dice_coef = float(len(files_a_sliced_b)) / (len(files_a) + len(files_b))
147          if (dice_coef > coef_threshold):
148              dice_stemmed_word_data.append((stemmed_word, word_a, word_b, dice_coef))
149
150          # Using MIM coef
151          mim_coef = float(len(files_a_sliced_b)) / (len(files_a) * len(files_b))
152          if (mim_coef > coef_threshold):
153              mim_stemmed_word_data.append((stemmed_word, word_a, word_b, mim_coef))
154
155          # Using EMIM coef
156          if len(files_a_sliced_b) > 0:
157              emim_coef = len(files_a_sliced_b) * \
158                          math.log(len(file_words_index) * len(files_a_sliced_b) /
       float(len(files_a) * len(files_b)))
159          else:
160              emim_coef = 0.0
161
162          if (emim_coef > coef_threshold):
163              emim_stemmed_word_data.append((stemmed_word, word_a, word_b, emim_coef))
164
165          # Chi-square
166          chi_sqr_coef = math.pow((len(files_a_sliced_b) - (float(len(files_a) * len(
       files_b)) / len(file_words_index))), 2) / \
```

```
167                          float ( len ( files _ a )  ∗  len ( files _ b ) )
168           if  ( chi _ sqr _ coef  >  coef _ threshold ) :
169               chi _ sqr _ stemmed _ word _ data . append (( stemmed _ word ,  word _ a ,  word _ b ,
       chi _ sqr _ coef ) )
170
171       if  len ( bigrams )  >  0 :  counter+=1
172       if  counter  >=  10 :  break
173
174
175
176 # Create  graph
177 dice _ stemmed _ word _ clusters  =  create _ graph ( dice _ stemmed _ word _ data ,  ' dice ' )
178 mim _ stemmed _ word _ clusters  =  create _ graph ( mim _ stemmed _ word _ data ,  ' mim ' )
179 emim _ stemmed _ word _ clusters  =  create _ graph ( emim _ stemmed _ word _ data ,  ' emim ' )
180 chi _ sqr _ stemmed _ word _ clusters  =  create _ graph ( chi _ sqr _ stemmed _ word _ data ,  ' chi _ sqr ' )
181
182 # Print  clusters
183 print _ stem _ cluster ( dice _ stemmed _ word _ clusters ,  ' dice ' )
184 print _ stem _ cluster ( mim _ stemmed _ word _ clusters ,  ' mim ' )
185 print _ stem _ cluster ( emim _ stemmed _ word _ clusters ,  ' emim ' )
186 print _ stem _ cluster ( chi _ sqr _ stemmed _ word _ clusters ,  ' chi _ sqr ' )
```

Listing 5: Python code for MLN1

I create the cluster using the graphs and extract the connected components from the graph. Figure 5 shows the example of the graph for the word 'feature'.
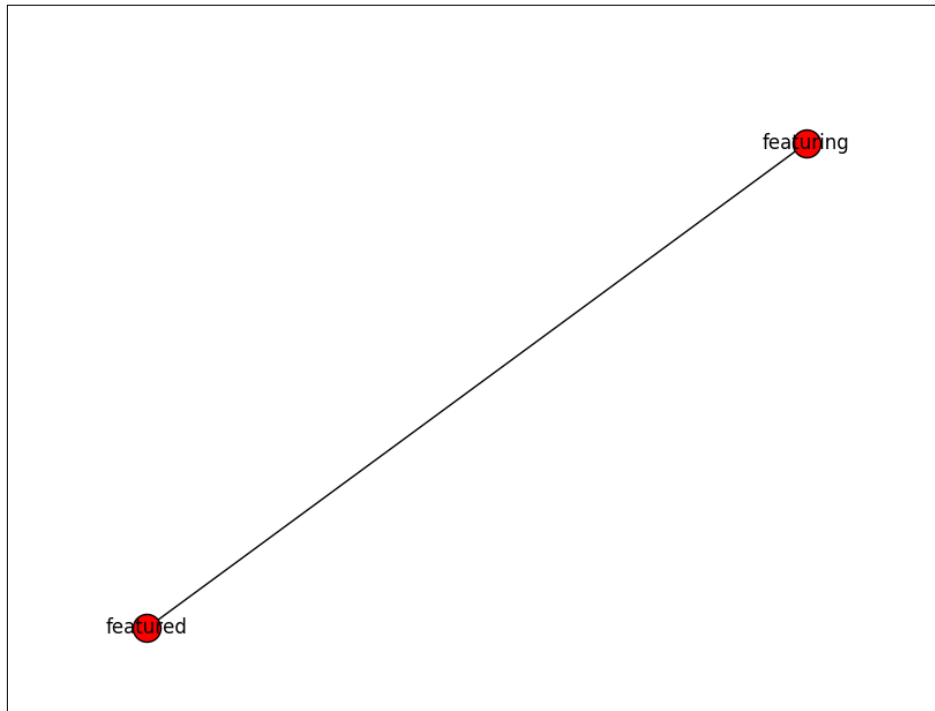


Figure 5: Connected component for the word 'feature'

From the output that I got, I see that difference association measurement return different clustering result. I also noticed that several terms do not have any other components connected to it. Figure 6 shows the cluster resulted from using Dice's Coefficient, MIM, EMIM, and Chi-square as the association measures.

16

```
Clusters using metric dice
==================================
provide : provide, provides
fin : fins, fin
feature : featured, featuring


Clusters using metric mim
==================================
provide : provide, provides
fin : fins, fin
feature : featured, featuring


Clusters using metric emim
==================================
provide : provide, provides
fin : fins, fin
feature : featured, featuring


Clusters using metric chi_sqr
==================================
daughter  : daughter, daughters
army   : armies, army
provide : provide, provides
type   : type, types
feature : featured, featuring
free   : freeing, free
standard  : standard, standards
street   : street, streets
fin : fins, fin
unit  : units, unit
```

Figure 6: The clustering result using Dice's Coefficient, MIM, EMIM, and Chi-square

## 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

Listing 6 shows the code to find the most strongly associated words for all 10 words that I have chosen.

```python
#!/usr/bin/python

import json
import os
import random
import html2text
import io
import math
import nltk

from krovetzstemmer import Stemmer as KrovetzStemmer
from nltk.corpus import stopwords


class AssociationMeasure:
    _document_words_index = {}
    _window_words_index_cache = {}
    _word_windows_index_cache = {}

    def __init__(self, dir_name):
        '''
        List all files in dir_name, and create document-words index
        :param dir_name:
        '''

        self._document_words_index = json.loads(io.open(dir_name).read())

    def get_window_words_index(self, window_size=None):
        '''
        Get window-words index
        :param window_size: number of words in a window, None or False if use
    document
        :return:
        '''

        # Re-build word_windows_index if window_size is not exists in
    _word_windows_index_cache
        if window_size not in self._window_words_index_cache:
            if not window_size:
                window_words_index = self._document_words_index
            else:
                window_words_index = {}
                for document, words in self._document_words_index.items():
                    for w in range(1, int(len(words) / window_size)):
                        start = (w - 1) * window_size
                        end = w * window_size
                        if end > len(words): end = len(words)

                        window_words = words[start:end]
                        window_words_index[document + '_part_' + str(w)] =
    window_words

            # Cache word_windows_index for each window_size
            self._window_words_index_cache[window_size] = window_words_index
```

```python
53
54
55            return self._window_words_index_cache[window_size]
56
57        def get_all_words(self, window_words_index):
58            '''
59            Get all words in a window−words index
60            :param window_words_index:
61            :return:
62            '''
63
64            all_words = []
65            for window, words in window_words_index.items():
66                all_words += words
67
68            all_words = set(all_words)
69            return all_words
70
71        def get_stemmed_words_index(self, window_words_index):
72            '''
73            Get stemmed−words index from window−words index
74            :param window_words_index:
75            :return:
76            '''
77
78            all_words = self.get_all_words(window_words_index)
79            stem_words_index = {}
80
81            krovetz = KrovetzStemmer()
82            for word in all_words:
83                # Stem word using krovetz
84                stemmed_word = krovetz.stem(word)
85
86                # Group by stemmed word
87                stem_words_index.setdefault(stemmed_word, [])
88                stem_words_index[stemmed_word].append(word)
89
90            return stem_words_index
91
92        def get_word_windows_index(self, window_size=None):
93            '''
94            Get word−windows index
95            :param window_size: number of words in a window, None or False if use
        document
96            :return:
97            '''
98
99            return json.loads(io.open('../6_1/1_2_index.txt').read())
100
101            # # Re−build word_windows_index if window_size is not exists in
        _word_windows_index_cache
102            # if window_size not in self._word_windows_index_cache:
103            #     window_words_index = self.get_window_words_index(window_size)
104            #     all_words = self.get_all_words(window_words_index)
105            #
106            #     word_windows_index = {}
107            #     for idx, word in enumerate(all_words):
108            #         print('{} of {}. Indexing word {}'.format(idx + 1, len(all_words),
        word.encode('utf−8')))
```

```python
109            #                print('=' * 30)
110            #
111            #                files = []
112            #                for file, words in window_words_index.items():
113            #                    if word in words:
114            #                        files.append(file)
115            #                word_windows_index[word] = sorted(set(files))
116            #
117            #
118            #        # Cache word_windows_index for each window_size
119            #        self._word_windows_index_cache[window_size] = word_windows_index
120            #
121            #
122            # return self._word_windows_index_cache[window_size]

124    def dice(self, word_a, word_b, window_size=None):
125        '''
126        Use dice coefficient
127        :param word_1:
128        :param word_2:
129        :param window_size:
130        :return:
131        '''
132
133        # Create window-word index
134        word_windows_index = self.get_word_windows_index(window_size)
135
136        # Lookup filename in word_files_index
137        files_a = word_windows_index[word_a]
138        files_b = word_windows_index[word_b]
139        files_a_sliced_b = list(set(files_b) & set(files_a))
140
141        dice_coef = float(len(files_a_sliced_b)) / (len(files_a) + len(files_b))
142        return dice_coef

144    def mim(self, word_a, word_b, window_size=None):
145        '''
146        User Mutual Information Measure
147        :param word_a:
148        :param word_b:
149        :param window_size:
150        :return:
151        '''
152
153        # Create window-word index
154        word_windows_index = self.get_word_windows_index(window_size)
155
156        # Lookup filename in word_files_index
157        files_a = word_windows_index[word_a]
158        files_b = word_windows_index[word_b]
159        files_a_sliced_b = list(set(files_b) & set(files_a))
160
161        mim_coef = float(len(files_a_sliced_b)) / (len(files_a) * len(files_b))
162        return mim_coef

164    def emim(self, word_a, word_b, window_size=None):
165        '''
166        Use Expected Mutual Information Measure
167        :param word_a:
```

```
168            :param word_b:
169            :param window_size:
170            :return:
171            '''
172
173            # Get window-words and word-windows index
174            window_words_index = self.get_window_words_index(window_size)
175            word_windows_index = self.get_word_windows_index(window_size)
176
177            # Lookup filename in word_files_index
178            files_a = word_windows_index[word_a]
179            files_b = word_windows_index[word_b]
180            files_a_sliced_b = list(set(files_b) & set(files_a))
181
182            if len(files_a_sliced_b) > 0:
183                emim_coef = len(files_a_sliced_b) * \
184                            math.log(len(window_words_index) * len(files_a_sliced_b) /
        float(len(files_a) * len(files_b)))
185            else:
186                emim_coef = 0.0
187
188            return emim_coef
189
190        def chi_square(self, word_a, word_b, window_size=None):
191            '''
192            Chi-Square Measure
193            :param word_a:
194            :param word_b:
195            :param window_size:
196            :return:
197            '''
198
199            # Get window-words and word-windows index
200            window_words_index = self.get_window_words_index(window_size)
201            word_windows_index = self.get_word_windows_index(window_size)
202
203            # Lookup filename in word_files_index
204            files_a = word_windows_index[word_a]
205            files_b = word_windows_index[word_b]
206            files_a_sliced_b = list(set(files_b) & set(files_a))
207
208            chi_sqr_coef = math.pow((len(files_a_sliced_b) -
209                                    (float(len(files_a) * len(files_b)) / len(
        window_words_index))), 2) / \
210                            float(len(files_a) * len(files_b))
211            return chi_sqr_coef
212
213
214 def get_most_associates(most_associated):
215     most_associated_sorted = {}
216     for word, associates in most_associated.items():
217         associates = sorted(associates, key=lambda x: x[1], reverse=True)
218         if len(associates) > 10: associates = associates[:10]
219
220         most_associated_sorted[word] = associates
221
222
223     return most_associated_sorted
224
```

```python
225
226  if __name__ == '__main__':
227      selected_words = ['abolishes', 'access', 'accommodate', 'accredited', 'sky',
228                        'railroad', 'calendar', 'airplane', 'airplane', 'bicycle']
229
230      measure = AssociationMeasure('../6_1/1_2_file_words_index.txt')
231
232      # Get all words and create all possible bigrams
233      all_word = measure.get_word_windows_index().keys()
234      all_bigrams = bigrams = list(nltk.bigrams(all_word))
235
236      # Calculate coocurence
237      most_associated_all_d = {}
238      most_associated_all_m = {}
239      most_associated_all_e = {}
240      most_associated_all_c = {}
241
242      # most_associated_5_d = {}
243      # most_associated_5_m = {}
244      # most_associated_5_e = {}
245      # most_associated_5_c = {}
246
247      print all_word
248      for word in selected_words:
249          for word_a, word_b in all_bigrams:
250              if word == word_a or word == word_b:
251                  dice_window_size_all = measure.dice(word_a, word_b)
252                  mim_window_size_all = measure.mim(word_a, word_b)
253                  emim_window_size_all = measure.emim(word_a, word_b)
254                  chi_square_window_size_all = measure.chi_square(word_a, word_b)
255
256                  # Use window-size = 5
257                  dice_window_size_5 = measure.dice(word_a, word_b, 5)
258                  mim_window_size_5 = measure.mim(word_a, word_b, 5)
259                  emim_window_size_5 = measure.emim(word_a, word_b, 5)
260                  chi_square_window_size_5 = measure.chi_square(word_a, word_b, 5)
261
262                  most_associated_all_d.setdefault(word, [])
263                  most_associated_all_m.setdefault(word, [])
264                  most_associated_all_e.setdefault(word, [])
265                  most_associated_all_c.setdefault(word, [])
266
267                  # most_associated_5_d.setdefault(word, [])
268                  # most_associated_5_m.setdefault(word, [])
269                  # most_associated_5_e.setdefault(word, [])
270                  # most_associated_5_c.setdefault(word, [])
271
272                  most_associated_all_d[word].append((word_a if word_b == word else
       word_b, dice_window_size_all))
273                  most_associated_all_m[word].append((word_a if word_b == word else
       word_b, mim_window_size_all))
274                  most_associated_all_e[word].append((word_a if word_b == word else
       word_b, emim_window_size_all))
275                  most_associated_all_c[word].append((word_a if word_b == word else
       word_b, chi_square_window_size_all))
276
277                  # most_associated_5_d[word].append((word_a if word_b == word else
       word_b, dice_window_size_5))
```

```
278                        # most_associated_5_m[word].append((word_a if word_b == word else
      word_b, mim_window_size_5))
279                        # most_associated_5_e[word].append((word_a if word_b == word else
      word_b, emim_window_size_5))
280                        # most_associated_5_c[word].append((word_a if word_b == word else
      word_b, chi_square_window_size_5))
281

282
283          # Sort ascending and get top 10
284          most_associated_all_d = get_most_associates(most_associated_all_d)
285          most_associated_all_m = get_most_associates(most_associated_all_m)
286          most_associated_all_e = get_most_associates(most_associated_all_e)
287          most_associated_all_c = get_most_associates(most_associated_all_c)
288
289          # most_associated_5_d = get_most_associates(most_associated_5_d)
290          # most_associated_5_m = get_most_associates(most_associated_all_d)
291          # most_associated_5_e = get_most_associates(most_associated_5_e)
292          # most_associated_5_e = get_most_associates(most_associated_5_e)
293
294          print most_associated_all_d
295
296      # # a. Find 10 words randomly
      ========================================================================
297      # stemmed_words_index = measure.get_stemmed_words_index(measure.
      get_window_words_index())
298      #
299      # # List only non-stop-words
300      # nltk.download('stopwords')
301      # stemmed_words= [stemmed_word for stemmed_word, words in stemmed_words_index.
      items()
302      #                      if stemmed_word not in stopwords.words('english')]
303      #
304      # # Find 10 stemmed-words randomly
305      # stemmed_words = random.sample(stemmed_words, 10)
306      #
307      #
308      # # b. Calculate co-occurence measure with window_size = 5 and window_size = full
      -document ===
309      # for word in stemmed_words:
310      #       # get words in stemmed-word-class
311      #       words = stemmed_words_index[word]
312      #       # create bigrams for words
313      #       bigrams = list(nltk.bigrams(words))
314      #
315      #       for word_a, word_b in bigrams:
316      #           # Use window-size = all document
317      #           dice_window_size_all = measure.dice(word_a, word_b)
318      #           mim_window_size_all = measure.mim(word_a, word_b)
319      #           emim_window_size_all = measure.emim(word_a, word_b)
320      #           chi_square_window_size_all = measure.chi_square(word_a, word_b)
321      #
322      #           # Use window-size = 5
323      #           dice_window_size_5 = measure.dice(word_a, word_b, 5)
324      #           mim_window_size_5 = measure.mim(word_a, word_b, 5)
325      #           emim_window_size_5 = measure.emim(word_a, word_b, 5)
326      #           chi_square_window_size_5 = measure.chi_square(word_a, word_b, 5)
327      #
328      #           print(u'Word class {} : {} - {} :'.format(word, word_a, word_b))
329      #
```

```
330    #                print('\tUsing dice coefficient')
331    #                print('\t\tWindow size = all ---> {}'.format(dice_window_size_all))
332    #                print('\t\tWindow size = {} ---> {}'.format(5, dice_window_size_5))
333    #
334    #                print('\tUsing MIM coefficient')
335    #                print('\t\tWindow size = all ---> {}'.format(mim_window_size_all))
336    #                print('\t\tWindow size = {} ---> {}'.format(5, mim_window_size_5))
337    #
338    #                print('\tUsing EMIM coefficient')
339    #                print('\t\tWindow size = all ---> {}'.format(emim_window_size_all))
340    #                print('\t\tWindow size = {} ---> {}'.format(5, emim_window_size_5))
341    #
342    #                print('\tUsing Chi-square coefficient')
343    #                print('\t\tWindow size = all ---> {}'.format(chi_square_window_size_all)
       )
344    #                print('\t\tWindow size = {} ---> {}'.format(5, chi_square_window_size_5)
       )
```

Listing 6: Association Measure

Tables below show the strong associated words for the chosen words.

Table 3: Strong associated words for 'school'

| MIM | EMIM | Chi Square | Dice |
|---|---|---|---|
| renaissance | poets | new | renaissance |
| sons | renaissance | poets | sons |
| formalism | new | group | poets |
| others | group | renaisaance | formalism |
| poets | son | southern | others |
| group | others | others | group |
| agrarians | formalism | sons | agrarians |
| artists | southern | formalism | artists |
| club | agrarians | club | club |
| southern | artists | agrarians | southern |

Table 4: Strong associated words for 'city'

| MIM | EMIM | Chi Square | Dice |
|---|---|---|---|
| tackle | back | back | tackle |
| area | tackle | new | back |
| founded | new | green | founded |
| county | green | bay | green |
| shopping | bay | end | bay |
| metropolitan | rams | packers | rams |
| privileges | packers | edit | packers |
| back | end | events | county |
| green | area | founded | new |

| Dice | MIM | EMIM | Chi-square |
|---|---|---|---|
| division | oxford | united | import |
| army | army | division | edit |
| oxford | divisions | army | united |
| united | division | oxford | id |
| th | birth | th | th |
| university | transportation | university | class |
| sovereign | allegiance | sovereign | division |
| press | sovereign | press | free |
| divisions | press | divisions | army |
| disambiguation | united | disambiguation | events |

Table 5: Strong associated words for 'states'

| Dice | MIM | EMIM | Chi-square |
|---|---|---|---|
| zero | zero | zero | zero |
| use | element | poetry | poetry |
| programming | programming | poets | poets |
| poets | systems | element | value |
| systems | array | may | also |
| would | however | would | index |
| however | would | also | use |
| array | wishing | programming | would |
| array | poets | however | would |
| element | use | use | may |

Table 6: Strong associated words for 'language'

| Dice | MIM | EMIM | Chi-square |
|---|---|---|---|
| zero | zero | zero | id |
| count | empty | count | zero |
| parser | set | parser | count |
| categories | defined | categories | main |
| expensive | shopping | expensive | value |
| bytes | comet | bytes | parser |
| poetry | particular | empty | bytese |
| size | zeros | set | categories |
| set | cardinality | poetry | expensive |
| asteroid | constant | defined | see |

Table 7: Strong associated words for 'function'

| Dice | MIM | EMIM | Chi-square |
|---|---|---|---|
| tackle | tackle | back | back |
| back | regiment | tackle | lions |
| regiment | soviet | lions | redskins |
| lions | artillery | redskinsd | giants |
| redskins | back | mm | guard |
| mm | lions | regiment | id |
| soviet | volunteer | giants | mm |
| infantry | redskins | guard | new |
| artillery | allegiance | infantry | tackle |
| giants | mount | soviet | infantry |

Table 8: Strong associated words for 'union'

| Dice | MIM | EMIM | Chi-square |
|---|---|---|---|
| limit | howitzer | limit | edit |
| preprocessor | song | preprocessor | count |
| node | area | node | limit |
| count | fellowships | count | preprocessor |
| kingdom | preprocessor | howitzer | node |
| championships | limit | kingdom | see |
| post | node | song | article |
| howitzer | gunfighter | championships | th |
| expanding | workshop | area | also |
| seasons | fas | expanding | league |

Table 9: Strong associated words for 'report'

| Dice | MIM | EMIM | Chi-square |
|---|---|---|---|
| zero | zero | zero | zero |
| buy | buy | buy | buy |
| limited | limited | limited | limited |
| exponentiation | exponentiation | exponentiation | exponentiation |
| four | four | four | four |
| sleep | sleep | sleep | sleep |
| oldest | oldest | oldest | oldest |
| whose | whose | whose | whose |
| calculate | calculate | calculate | calculate |
| segments | segments | segments | segments |

Table 10: Strong associated words for 'fish'

| Dice | MIM | EMIM | Chi-square |
|---|---|---|---|
| singles | singles | singles | edit |
| championships | elections | championships | championships |
| elections | procterate | elections | round |
| protectorate | permitted | act | main |
| act | parliamentary | protectorate | article |
| candidadtes | stand | candidates | singles |
| allowing | championships | allowing | act |
| permitted | candidates | permitted | defeated |
| defeated | allowing | defeated | elections |
| parliamentary | act | parliamentary | candidates |

Table 11: Strong associated words for 'women'

# 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.