

Assignment 3

CS834 Introduction to Information Retrieval

Fall 2017

Polykarpou Thomadakis

November 8, 2017

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 (Dices 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 this assignment, I used the code I used in assignment 2 to produce the inverted index. Then, I sort the index alphabetically and pick the first 1000 words and create stems from each one of them using the krovetz stemmer. Once I have created the stems, I create bigrams from them in order to find co-occurrences in the given collection and apply Dice's coefficient to create the stemming clusters. I set the threshold to 0.01 and then extracted the connected subgraphs produced by the code. The code used can be seen in listing 1.

```
1 import os
# import os.path
3 import csv
from bs4 import BeautifulSoup
5 import nltk
import argparse
7 import krovetzstemmer
9 def hasNumbers(inputString):
    return any(char.isdigit() for char in inputString)
11 tokenizer = nltk.RegexpTokenizer(r '\w+')
13 parser = argparse.ArgumentParser('invindex')
15 group = parser.add_mutually_exclusive_group(required= True)
```

```

group.add_argument(' -l ', '--list', nargs='+', help='List of documents to parse',
    required=False)
17 group.add_argument(' -p ', '--path', help='Path to directory to parse all of its files',
    required=False)
args = parser.parse_args();
19
20 html_files = []
21 if(args.path):
22     root_dir = args.path
23     for root, dirs, files in os.walk(root_dir):
24         for file in files:
25             if file.endswith('.html'):
26                 html_files.append(os.path.join(root, file))
27 else:
28     html_files = [x for x in args.list if x.endswith('.html')]
29
30 inv_index = {}
31
32 if not os.path.exists("./inv_file.csv"):
33
34     for file in html_files:
35         print 'Indexing file: '+file
36         with open(file) as f:
37             soup = BeautifulSoup(f.read(), 'html.parser')
38             tokens= tokenizer.tokenize(soup.get_text())
39             for token in tokens:
40                 token = token.encode("utf-8")
41                 if token not in inv_index:
42                     inv_index[token] = set()
43                     inv_index[token].add(os.path.split(file)[1])
44
45         outfile = open("inv_file.csv", "w")
46         for token in sorted(inv_index.iterkeys()):
47             if hasNumbers(token) or token.isupper():
48                 continue;
49             outfile.write(token+",")
50             i=0
51             outfile.write(',')
52             for file in inv_index[token]:
53                 i+=1
54                 outfile.write(file)
55                 if(i != len(inv_index[token])):
56                     outfile.write(",")
57             outfile.write("\n")
58
59 krovetz = krovetzstemmer.Stemmer()
60
61 stem_words = {}
62 index = {}
63 count = 0;
64 with open('inv_file.csv') as csvfile:
65     for c in csvfile:
66         word, files = c.split(",",1)
67         files = files.replace("\r",",")
68         files = files.strip()
69         index[word] = files.split(",")
70         if count == 1000:
71             continue;
72         stem = krovetz.stem(word)
73         if stem not in stem_words:
74             stem_words[stem] = []
75             stem_words[stem].append(word)
76         count+=1
77
78 coef_thres = 0.001
79 before_dice = []

```

```

81 after_dice = []
82 for stem in sorted(stem_words):
83     bigrams = nltk.bigrams(stem_words[stem])
84     for a,b in bigrams:
85         f_a = index[a]
86         len_f_a = len(f_a)
87         f_b = index[b]
88         len_f_b = len(f_b)
89         f_ab = set(index[a]) & set(index[b])
90         len_f_ab = len(f_ab)
91         before_dice.append((stem,a,b))
92         coef = float(2*len_f_ab)/(len_f_a+len_f_b)
93         if coef > coef_thres:
94             if stem not in after_dice:
95                 after_dice[stem] = set()
96                 after_dice[stem].add(a)
97                 after_dice[stem].add(b)
98
99 for stem in after_dice:
100    print stem +": {",
101    for word in after_dice[stem]:
102        print word,
103    print "}"

```

Listing 1: Script to extract the stem clusters using Dice's coefficient

The clusters created by the algorithm are presented in table 1.

Table 1: Clusters created using the Dice's coefficient on the small wikipedia collection

academician	Academician,Academicians
abbott	Abbotts,Abbott
abstraction	Abstraction,Abstractions
abugida	Abugidas,Abugida
acquisition	Acquisitions,Acquisition
ace	Aces,Acer
actor	Actors,Actor
acronym	Acronym,Acronyms
access	Access,Accessed
account	Accounting,Account,Accounts
abbey	Abbeys,Abbey
abbreviation	Abbreviation,Abbreviations
accident	Accident,Accidents
abt	Abts,Abt
absurd	Absurdity,Absurd
acre	Acres,Acre
atpase	ATPase,ATPases
activity	Activities,Activity
achievement	Achievements,Achievement
adaboost	AdaBoost,Adaboost

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

For the edit distance calculator I used the python code found in <http://norvig.com/spell-correct.html>, specifically edits1 and edits2, that which given a misspelled word apply 1 and 2 letter transformations to the misspelled word with the purpose of generating a valid one. Once I collect the set of words generated this way, I keep only those that are valid words by checking if they exist in a big collection of words in file *big.txt* found in the repository (adopted from the same site noted above). Once the valid words have been collected, I compute the probability the frequency that each of those words appears, and finally correct the misspelled word with the most frequent word of the set of valid words collected. The code for this assignment appears in listing 2.

```

1 import re
2 import sys
3 from collections import Counter
4
5 text = open('big.txt').read()
6 words = re.findall(r'\w+', text.lower())
7 word_freq = Counter(words)
8 num_words = 0
9
10 if len(sys.argv) != 2:
11     print 'Please input word to correct'
12     exit()
13
14 input_word = sys.argv[1].lower()
15
16 for counter in word_freq.values():
17     num_words += counter
18
19 def edits1(word):
20     """All edits that are one edit away from 'word'."""
21     letters = 'abcdefghijklmnopqrstuvwxyz'
22     splits = [(word[:i], word[i:]) for i in range(len(word) + 1)]
23     deletes = [L + R[1:] for L, R in splits if R]
24     transposes = [L + R[1] + R[0] + R[2:] for L, R in splits if len(R) > 1]
25     replaces = [L + c + R[1:] for L, R in splits if R for c in letters]
26     inserts = [L + c + R for L, R in splits for c in letters]
27     return set(deletes + transposes + replaces + inserts)
28
29 def edits2(word):
30     """All edits that are two edits away from 'word'."""
31     return (e2 for e1 in edits1(word) for e2 in edits1(e1))
32
33 def probability(word):
34     return float(word_freq[word]) / num_words
35
36
37 correct_words = set()
38
39 for x in edits1(input_word) or edits2(input_word):
40     if x in word_freq:
41         correct_words.add(x)
42
43 max_ = 0
44 corrected_word = input_word
45 for p in correct_words:
46     if word_freq[p] > max_:
47         max_ = word_freq[p]
48         corrected_word = p
49

```

```
print corrected_word
```

Listing 2: Script to correct a misspelled word based on the noisy channel model

Figure 1 demonstrates the spell correction program listed above.

```
pthomadakis@e-3100-ubuntu:~/Desktop/CS834/A3$ python spell_correction.py lodon
london
pthomadakis@e-3100-ubuntu:~/Desktop/CS834/A3$ python spell_correction.py virgninia
virginia
pthomadakis@e-3100-ubuntu:~/Desktop/CS834/A3$ python spell_correction.py wasington
washington
pthomadakis@e-3100-ubuntu:~/Desktop/CS834/A3$ python spell_correction.py cicago
chicago
pthomadakis@e-3100-ubuntu:~/Desktop/CS834/A3$ python spell_correction.py milwauke
milwaukee
```

Figure 1: Sample output of the spell correction script

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

The code for the Galago search engine can be found in <https://storage.googleapis.com/google-code-archive-downloads/v2/code.google.com/galagosearch/galagosearch-1.04-src.tar.gz>. The code for the snippet generation resides in the file: *galagosearch-1.04/galagosearch-core/src/main/java/org/galagosearch/core/store/SnippetGenerator.java*. The algorithm can be described as the following 6 steps:

1. Tokenize the document into terms and maintain them and their positions.
2. Find all the terms in the document that match the ones of the query.
3. When a term is matched, retrieve a small region of the document that contains the term.
4. Check for overlapping regions and merge them if any.
5. Keep adding document regions to the snippet, until the maximum size is reached.
6. Combine all regions to a single html string.

Galago includes snippets only when the exact term is found in the document. Thus, when a small document is indexed there is a lower probability that a term will appear and as a result a lower probability that a snippet will be generated. Also, Galago does not seem to use stemming, meaning that it is even less probable for a snippet to be generated. In my opinion, stemming would increase the quality of snippet generation as it does for searching. If this also fails, it would be a good attempt to just display the first few sentences of the document. Since it has already been marked as relevant, it makes sense that at the introduction of the document there will be some reference relevant to the query term.

Question MLN1

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

Answer

This assignment can be answered directly using the code from question 6.1. To add more value to the code used there I added the stem clusters produced when other association metrics are used, namely Dice's coefficient, EMIM, MIM and Chi-square.

The ten words I randomly chose are:

“academician”, “university”, “mind”, “synonym”, “monitor”, “owe”, “top”, “charge”, “hit”, “diverge”. Figure 2 shows the classes created from each metric using threshold 0.01. The code used for this assignment is shown in listing 3

```
import os
# import os.path
import csv
from bs4 import BeautifulSoup
import nltk
import argparse
import krovetzstemmer
import math
import random

chosen_words = [ "academician" , "university" , "mind" , "synonym" , "monitor" , "owe" ,
    "top" , "charge" , "hit" , "diverge" ]
def hasNumbers(inputString):
    return any(char.isdigit() for char in inputString)

def printIndex(index):
    for stem in chosen_words:
        print stem +": {",
        if( stem in index):
            for word in index[stem]:
                print word,
            print "}"
        else:
            print "}"

def dice(a,b):
    f_a = index[a]
    len_f_a = len(f_a)
    f_b = index[b]
    len_f_b = len(f_b)
    f_ab = set(index[a]) & set (index[b])
    len_f_ab = len(f_ab)
    # before_dice.append((stem,a,b))
    coef = float(len_f_ab)/(len_f_a+len_f_b)
    if coef > coef_thres:
        return True
    return False

def mim(a,b):
    f_a = index[a]
    len_f_a = len(f_a)
    f_b = index[b]
    len_f_b = len(f_b)
    f_ab = set(index[a]) & set (index[b])
    len_f_ab = len(f_ab)
    coef = float(len_f_ab)/(len_f_a*len_f_b)
    if coef > coef_thres:
        return True
    return False

def emim(a,b,N):
    f_a = index[a]
    len_f_a = len(f_a)
    f_b = index[b]
    len_f_b = len(f_b)
```

```

56     f_ab = set(index[a]) & set (index[b])
57     len_f_ab = len(f_ab)
58     if len_f_ab == 0:
59         return False
60     coef = len_f_ab*math.log(N*(float(len_f_ab)/(len_f_a*len_f_b)))
61     if coef > coef_thres:
62         return True
63     return False

64 def chi_sqr(a,b,N):
65     f_a = index[a]
66     len_f_a = len(f_a)
67     f_b = index[b]
68     len_f_b = len(f_b)
69     f_ab = set(index[a]) & set (index[b])
70     len_f_ab = len(f_ab)
71     coef = len_f_ab-1/float(N)*len_f_a*len_f_b
72     coef = coef*coef/(len_f_a*len_f_b)
73     if coef > coef_thres:
74         return True
75     return False

76 tokenizer = nltk.RegexpTokenizer(r'\w+')
77
78 parser = argparse.ArgumentParser('invindex')
79 group = parser.add_mutually_exclusive_group(required= True)
80 group.add_argument('-l', '--list', nargs='+', help='List of documents to parse',
81     required=False)
82 group.add_argument('-p', '--path', help='Path to directory to parse all of its files',
83     required=False)
84 args = parser.parse_args();

85 html_files = []
86 if(args.path):
87     root_dir = args.path
88     for root,dirs,files in os.walk(root_dir):
89         for file in files:
90             if file.endswith('.html'):
91                 html_files.append(os.path.join(root, file))
92 else:
93     html_files = [x for x in args.list if x.endswith('.html')]
94
95 inv_index = {}

96 if not os.path.exists("./inv_file.csv"):
97
98     for file in html_files:
99         print 'Indexing file: '+file
100        with open(file) as f:
101            soup = BeautifulSoup(f.read(), 'html.parser')
102            tokens= tokenizer.tokenize(soup.get_text())
103            for token in tokens:
104                token = token.encode("utf-8")
105                if token not in inv_index:
106                    inv_index[token] = set()
107                    inv_index[token].add(os.path.split(file)[1])

108
109        outfile = open("inv_file.csv","w")
110        for token in sorted(inv_index.iterkeys()):
111            if hasNumbers(token) or token.isupper():
112                continue;
113            outfile.write(token+",")
114            i=0
115            outfile.write(",")
116            for file in inv_index[token]:
117                i+=1

```

```

120     outfile.write(file)
121     if(i != len(inv_index[token])):
122         outfile.write(",")
123         outfile.write("\n")
124
125 krovetz = krovetzstemmer.Stemmer()
126
127 stem_words = {}
128 index = {}
129 count = 0;
130 with open('inv_file.csv') as csvfile:
131     for c in csvfile:
132         word, files = c.split(",",1)
133         files = files.replace("\r\n","")
134         files = files.strip()
135         index[word] = files.split(",")
136         # if count == 1000:
137             # continue;
138         stem = krovetz.stem(word)
139         if stem not in stem_words:
140             stem_words[stem] = []
141             stem_words[stem].append(word)
142             # count+=1
143
144 # print chosen_words
145 coef_thres = 0.01
146 dice_index = {}
147 mim_index = {}
148 emim_index = {}
149 chi_sqr_index = {}
150 for stem in sorted(stem_words):
151     bigrams = nltk.bigrams(stem_words[stem])
152     for a,b in bigrams:
153         if dice(a,b):
154             if stem not in dice_index:
155                 dice_index[stem] = set()
156                 dice_index[stem].add(a)
157                 dice_index[stem].add(b)
158             if mim(a,b):
159                 if stem not in mim_index:
160                     mim_index[stem] = set()
161                     mim_index[stem].add(a)
162                     mim_index[stem].add(b)
163             if emim(a,b,len(html_files)):
164                 if stem not in emim_index:
165                     emim_index[stem] = set()
166                     emim_index[stem].add(a)
167                     emim_index[stem].add(b)
168             if chi_sqr(a,b,len(html_files)):
169                 if stem not in chi_sqr_index:
170                     chi_sqr_index[stem] = set()
171                     chi_sqr_index[stem].add(a)
172                     chi_sqr_index[stem].add(b)
173
174 # print "N = ",len(html_files)
175 print "Dice = "
176 printIndex(dice_index)
177 print "MIM = "
178 printIndex(mim_index)
179 print "EMIM = "
180 printIndex(emim_index)
181 print "CHI_SQR = "
182 printIndex(chi_sqr_index)

```

Listing 3: Script to generate stem classes for the 10 chosen words

```

=====
Dice =====
academician: { Academician academician academicians Academicians }
university: { Universities universities University university }
mind: { mind minds Minds }
synonym: { synonymize synonym synonymized }
monitor: { monitored monitoring Monitoring monitor monitors }
owe: { owed owes }
top: { top topped }
charge: { Charges charges Charge charge charging }
hit: { hitting Hits hit hitter hits Hitting hitters Hit }
diverge: { diverging diverges }
=====
MIM =====
academician: { Academician academician academicians Academicians }
university: { }
mind: { }
synonym: { synonymize synonym synonymized }
monitor: { }
owe: { }
top: { }
charge: { Charges Charge }
hit: { hitters hitter }
diverge: { diverging diverges }
=====
EMIM =====
academician: { Academician academician academicians Academicians }
university: { Universities universities University university }
mind: { mind minds Minds }
synonym: { synonymize synonym synonymized }
monitor: { monitored monitoring Monitoring monitor monitors }
owe: { owed owes }
top: { top topped Topped }
charge: { Charges charge charging charges Charge Charging }
hit: { hitting Hits hit hitter hits Hitting hitters Hit }
diverge: { diverging diverges }
=====
CHI_SQR =====
academician: { Academician academician academicians Academicians }
university: { Universities universities University university }
mind: { mind minds }
synonym: { synonymize synonym synonymized }
monitor: { monitoring monitored }
owe: { }
top: { }
charge: { Charges charges Charge charge }
hit: { Hits hit hitter hits hitters Hit }
diverge: { diverging diverges }

```

Figure 2: Stem classes produced for different association metrics

Question MLN2

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

Answer

The script that was written for this assignment can be found in listing 4. It uses the inverted index of the previous assignments (or regenerates it if it does not exist) to calculate the association measures using the metrics: MIM, EMIM, Chi-squared and Dice's coefficient. The ten words I randomly chose are: "instrument", "university", "book", "America", "fish", "dinosaur", "car", "food", "hit", "soccer".

```
1 import os
2 import csv
3 from bs4 import BeautifulSoup
4 import nltk
5 import argparse
6 import krovetzstemmer
7 import math

9 chosen_words = [ "instrument" , "university" , "book" , "America" , "fish" , "dinosaur" , "car" , "food" , "hit" , "soccer" ]

11 def hasNumbers(inputString):
12     return any(char.isdigit() for char in inputString)

13 def sortByValue(tup):
14     return tup[1]

17 def insertAssociate(a,b,metric,coef):
18     if a not in associates[metric]:
19         associates[metric][a] = []
20     associates[metric][a].append([b,coef])

21 def printIndex(index):
22     for stem in chosen_words:
23         print stem +": {" ,
24         if( stem in index):
25             for word in index[stem]:
26                 print word ,
27                 print "}"
28         else:
29             print "}"

31 def dice(a,b):
32     f_a = index[a]
33     len_f_a = len(f_a)
34     f_b = index[b]
35     len_f_b = len(f_b)
36     f_ab = set(index[a]) & set(index[b])
37     len_f_ab = len(f_ab)
38     coef = float(len_f_ab)/(len_f_a+len_f_b)
39     insertAssociate(a,b,'dice',coef)

41 def mim(a,b):
42     f_a = index[a]
43     len_f_a = len(f_a)
44     f_b = index[b]
45     len_f_b = len(f_b)
46     f_ab = set(index[a]) & set(index[b])
47     len_f_ab = len(f_ab)
48     coef = float(len_f_ab)/(len_f_a*len_f_b)
49     insertAssociate(a,b,'mim',coef)

51
```

```

53     def emim(a,b,N):
54         f_a = index[a]
55         len_f_a = len(f_a)
56         f_b = index[b]
57         len_f_b = len(f_b)
58         f_ab = set(index[a]) & set(index[b])
59         len_f_ab = len(f_ab)
60         if len_f_ab == 0:
61             coef = 0.0
62             insertAssociate(a,b,'emim',coef)
63             return
64         coef = len_f_ab*math.log(N*(float(len_f_ab)/(len_f_a*len_f_b)))
65         insertAssociate(a,b,'emim',coef)
66
67     def chi_sqr(a,b,N):
68         f_a = index[a]
69         len_f_a = len(f_a)
70         f_b = index[b]
71         len_f_b = len(f_b)
72         f_ab = set(index[a]) & set(index[b])
73         len_f_ab = len(f_ab)
74         coef = len_f_ab-1/float(N)*len_f_a*len_f_b
75         coef = coef*coef/(len_f_a*len_f_b)
76         insertAssociate(a,b,'chi',coef)
77
78     tokenizer = nltk.RegexpTokenizer(r'\w+')
79
80     parser = argparse.ArgumentParser('invindex')
81     group = parser.add_mutually_exclusive_group(required= True)
82     group.add_argument('-l','--list', nargs='+', help='List of documents to parse',
83                         required=False)
84     group.add_argument('-p','--path', help='Path to directory to parse all of its files',
85                         required=False)
86     args = parser.parse_args();
87
88     html_files = []
89     if(args.path):
90         root_dir = args.path
91         for root,dirs,files in os.walk(root_dir):
92             for file in files:
93                 if file.endswith('.html'):
94                     html_files.append(os.path.join(root, file))
95     else:
96         html_files = [x for x in args.list if x.endswith('.html')]
97
98     inv_index = {}
99
100    if not os.path.exists("./inv_file.csv"):
101        for file in html_files:
102            print 'Indexing file: '+file
103            with open(file) as f:
104                soup = BeautifulSoup(f.read(), 'html.parser')
105                tokens= tokenizer.tokenize(soup.get_text())
106                for token in tokens:
107                    token = token.encode("utf-8")
108                    if token not in inv_index:
109                        inv_index[token] = set()
110                        inv_index[token].add(os.path.split(file)[1])
111
112                outfile = open("inv_file.csv","w")
113                for token in sorted(inv_index.iterkeys()):
114                    if hasNumbers(token) or token.isupper():
115                        continue;
116                    outfile.write(token+",")
117                    i=0

```

```

117     outfile.write(',')
118     for file in inv_index[token]:
119         i+=1
120         outfile.write(file)
121         if(i != len(inv_index[token])):
122             outfile.write(",")
123             outfile.write("\n")
124
125     krovetz = krovetzstemmer.Stemmer()
126
127     stem_words = {}
128     index = {}
129     count = 0;
130     with open('inv_file.csv') as csvfile:
131         for c in csvfile:
132             word, files = c.split(",",1)
133             files = files.replace("\r\n","")
134             files = files.strip()
135             index[word] = files.split(",")
136
137     words = index.keys()
138
139     associates = {}
140     associates['dice'] = {}
141     associates['mim'] = {}
142     associates['emim'] = {}
143     associates['chi'] = {}
144
145     for choice in chosen_words:
146         for word in words:
147             if word != choice:
148                 dice(choice,word)
149                 mim(choice,word)
150                 emim(choice,word,len(html_files))
151                 chi_sqr(choice,word,len(html_files))
152
153     per_word_associates = {}
154     out = open("associates","w")
155     for choice in chosen_words:
156         per_word_associates[choice] = {}
157         per_word_associates[choice]['dice'] = []
158         per_word_associates[choice]['mim'] = []
159         per_word_associates[choice]['emim'] = []
160         per_word_associates[choice]['chi'] = []
161         per_word_associates[choice]['dice'] = sorted(associates['dice'][choice],key=
162             sortByValue,reverse=True)[:10]
163         per_word_associates[choice]['mim'] = sorted(associates['mim'][choice],key=
164             sortByValue,reverse=True)[:10]
165         per_word_associates[choice]['emim'] = sorted(associates['emim'][choice],key=
166             sortByValue,reverse=True)[:10]
167         per_word_associates[choice]['chi'] = sorted(associates['chi'][choice],key=
168             sortByValue,reverse=True)[:10]
169         out.write(choice+"\n")
170         out.write("Dice,MIM,EMIM,Chi\n")
171         for i in range(0,10):
172             out.write(per_word_associates[choice]['dice'][i][0]+","+
173             per_word_associates[choice]['mim'][i][0]+","+
174             per_word_associates[choice]['emim'][i][0]+","+
175             per_word_associates[choice]['chi'][i][0]+\n")

```

Listing 4: Script to generate top 10 associated words for the 10 chosen words

The resulting associated words are shown in tables 2-11. We can see that the associated words can differ significantly depending on the metric that is used to measure association. The most significant difference in the table is noticed on the MIM metric, where almost none of the words is common with the other ones. That is probably caused by the fact that this metric favors low-frequency terms,

since $n_a * n_b$ will increase very fast for frequent words a,b but their co-occurrence remains the same. As a result, the total value will be very low compared to low-frequency.

Table 2: Associated words for word “instrument”

instrument			
Dice	MIM	EMIM	Chi
instruments	Mortem	instruments	instruments
signature	Serwaczynski	musical	timbres
Perspectives	conductive	music	Buggles
Furrykef	Jbonneau	such	ney
scores	Desinicization	some	Postmodernism
Instrument	Arhielanto	used	disclosed
sensitive	Dentine	up	Ulmanor
musical	Artisan	heavily	halts
Face	Dannycali	different	Crazed
timbres	SatellitesHidden	various	Maxillofacial

Table 3: Associated words for word “university”

university			
Dice	MIM	EMIM	Chi
institution	Tootie	University	institution
universities	Pamplemousse	education	universities
campus	Mortem	college	education
academic	Velocisaurus	institution	campus
education	Cheiro	students	academic
college	Babineau	campus	college
student	Wingwangwo	academic	postgraduate
colleges	Malfatti	universities	Faculties
Engineering	majr	College	University
students	Thegreyanomaly	student	colleges

Table 4: Associated words for word “book”

book			
Dice	MIM	EMIM	Chi
published	Simka	published	published
books	Mortem	books	books
story	SparrowsWing	story	story
written	fiction	his	novels
novel	Pvasiliadis	written	novel
author	Mucks	he	comic
wrote	Norilana	that	author
Books	Grecque	novel	fiction
series	author	they	Publication
Book	equilateral	author	wrote

Table 5: Associated words for word “America”

America			
Dice	MIM	EMIM	Chi
North	Tootie	North	North
American	Mortes	American	Europe
Europe	Mortem	New	Asia
New	seamier	show	American
York	Sergeants	States	Oceania
through	Pinkerton	hide	Mexico
States	Mucko	var	Brazil
most	Zizaniopsis	tocShowText	York
South	Velocisaurus	showTocToggle	Africa
began	Amirthanayagam	tocHideText	New

Table 6: Associated words for word “fish”

fish			
Dice	MIM	EMIM	Chi
Actinopterygii	hukim	species	Actinopterygii
Fish	Keulemans	Actinopterygii	FishBase
Chordata	desiccant	water	Ranier
freshwater	Jabr	Chordata	Osteichthyes
species	Sennet	Fish	Fishes
Phylum	Percopsiformes	Phylum	Fish
fins	bellow	classification	Pauly
Animalia	arsēn	Animalia	fins
FishBase	Protista	Scientific	fin
Ranier	bobwhite	Family	Froese

Table 7: Associated words for word “dinosaur”

dinosaur			
Dice	MIM	EMIM	Chi
Dinosaurs	Velocisaurus	dinosaurs	Dinosaurs
dinosaurs	Cucumber	Dinosaurs	dinosaurs
theropod	Gurche	Firsfron	theropod
Dinosauria	jil	theropod	Dinosauria
Cretaceous	candelariensis	Dinosauria	Saurischia
Sauropsida	sauropodomorph	Spencer	Cretaceous
Saurischia	edwardsorum	Cretaceous	Megalosaurus
Superorder	Zanclodon	Sauropsida	Theropoda
Jurassic	Baricco	Dinosaur	Sauropsida
Dodson	Diplodocus	Superorder	Dodson

Table 8: Associated words for word “car”

car			
Dice	MIM	EMIM	Chi
cars	Mortem	into	cars
vehicle	seamier	cars	vehicle
Car	Sergeants	they	Car
driver	SparrowsWing	through	driver
driving	Hafid	their	driving
vehicles	Mwstorlie	such	vehicles
hotel	dna	been	wheels
road	equilateral	new	hotel
get	titanium	have	automobile
drive	Alor	two	Shoppes

Table 9: Associated words for word “food”

food			
Dice	MIM	EMIM	Chi
drink	Behaalotech	have	drink
water	water	they	consume
themselves	Sergeants	there	diet
eat	canes	than	foods
roughly	Food	their	diarrhea
Food	Hif	more	prey
source	eat	most	roughly
significant	hukim	well	Food
low	Sameerkale	them	eating
areas	Vernooykill	are	water

Table 10: Associated words for word “hit”

hit			
Dice	MIM	EMIM	Chi
song	comically	song	song
re	Haynos	re	re
behind	Valli	second	hits
debut	Sedaka	released	getting
next	Shockey	their	fans
getting	Mucks	his	behind
featured	Mwstorlie	single	next
single	Vipiteno	off	single
fans	Akroyd	next	Billboard
video	Alou	up	debut

Table 11: Associated words for word “soccer”

soccer			
Dice	MIM	EMIM	Chi
Football	Karaliopolous	football	Football
Goals	Hig	Football	football
football	Visvliet	league	Goals
Gls	Bobae	club	Gls
App	Estay	team	App
league	NellieMcKay	Club	league
footballers	Ide	Goals	footballers
domestic	Spudders	Cup	domestic
goals	Thosol	domestic	club
Playing	Seade	goals	goals