

山东大学 计算机科学与技术 学院

信息检索与数据挖掘 课程实验报告

学号:201600150109	姓名: 沈棋韬	班级: 16 人工智能班
实验题目: Inverted index and Boolean Retrieval Model		
<p>实验内容:</p> <p>Use 30548 tweets to create inverted index and realize Boolean Retrieval Model.</p>		
<p>实验过程中遇到和解决的问题:</p> <p>(记录实验过程中遇到的问题, 以及解决过程和实验结果。可以适当配以关键代码辅助说明, 但不要大段贴代码。)</p> <p>First we need to load these tweets into memory. Because it is stored with json, I use</p> <pre>for line in file:     tweets.append(json.loads(line))</pre> <p>to generate a list named tweets which contains a dictionary of every tweet.</p> <p>Then, the texts of these tweets need to be pre-processd to get rid of meaningless symbols and correct some type errors. After this, I create a list (doc_word) to store the words of every tweets.</p> <p>When the preparation is done, I can start to count the frequency of every word in these tweets.</p> <pre>for tweet in doc_word:     tweet_num+=1     for word in tweet:         if word in word_dic.keys():             #列表的第一个值为频数             word_dic[word][0]+=1             word_dic[word].append(tweet_num)         else:             word_dic[word]=[]             word_dic[word].append(1)             word_dic[word].append(tweet_num)</pre> <p>For every tweet in the list, I traverse every word in it and add 1 to the frequency and append the number of tweet if it already exists in the word_dic or create a Key with this word and append the number</p>		

of the tweet and set its frequency to 1 if it is not in the word\_dic.

The structure of doc\_word:

```
['how',  
 'to',  
 'make',  
 'money',  
 'from',  
 'horse',  
 'racing',  
 'with',  
 'betting',  
 'systems',  
 'http',  
 'bitlyhtpzz3'],  
 ['sales',  
 'rise',  
 'lifting',  
 'mcdonald',  
 'profit',  
 'new',  
 'menu',  
 'items',  
 'helped',  
 'earnings',  
 'grow',  
 'and',  
 'overcome',  
 'slow',  
 'decembers',  
 'http',  
 'bitlyfip3ot'],
```

The structure of word\_dic:

```
'between': [132,  
127,  
190,  
248,  
493,  
2212,  
2622,  
2743,  
2843,  
2955,  
2972,  
2994,  
3035,  
3442,  
3579,  
3583,  
3888,  
4192,  
5214,  
5224,  
5643,  
6524,  
6580,  
7343,  
7420,  
7600
```

Every word in these tweets is a key of the dictionary. The first number of the list is the frequency of the word and the other numbers are the tweets in which the word appears.

Merge function: analyze the query and then invoke one of three functions(*mergeAnd*, *mergeOr* and *mergeNot* respectively). For example, if the query looks like “between and home”, the merge function will invoke *mergeAnd* and “between” and “home” will be the two parameters of the function.

The definition of these functions are as follows:

```

62 def merge(query):
63     words=nlk.word_tokenize(query)
64     words=process(words)
65     if 'and' in query:
66         # w1=process(words[0])
67         w1=words[0]
68         # w2=process(words[2])
69         w2=words[2]
70         return mergeAnd(w1,w2)
71     elif 'or' in query:
72         # w1=process(words[0])
73         w1=words[0]
74         # w2=process(words[2])
75         w2=words[2]
76         return mergeOr(w1,w2)
77     elif 'not' in query:
78         # w1=process(words[1])
79         w1=words[1]
80         return mergeNot(w1)
81     else:
82         print('wrong query')
83
84
85 def mergeAnd(w1,w2):
86     if w1 not in word_bag or w2 not in word_bag:
87         print('No result')
88         return
89     lst=[]
90     p1=1
91     p2=1
92     while True:
93         if p1==len(word_dic[w1]):
94             break
95         if p2==len(word_dic[w2]):
96             break
97         if word_dic[w1][p1]<word_dic[w2][p2]:
98             p1+=1
99         elif word_dic[w1][p1]>word_dic[w2][p2]:
100             p2+=1
101         elif word_dic[w1][p1]==word_dic[w2][p2]:
102             lst.append(word_dic[w1][p1])
103             p1+=1
104             p2+=1
105         else:
106             pass
107     return lst

```

```

108 def mergeOr(w1,w2):
109     if w1 not in word_bag and w2 not in word_bag:
110         print('No result')
111         return
112     lst=[]
113     p1=1
114     p2=1
115     while True:
116         if word_dic[w1][p1]<word_dic[w2][p2]:
117             lst.append(word_dic[w1][p1])
118             p1+=1
119             if p1==len(word_dic[w1]):
120                 while p2<len(word_dic[w2]):
121                     lst.append(word_dic[w2][p2])
122                     p2+=1
123                 return lst
124         if word_dic[w1][p1]>word_dic[w2][p2]:
125             lst.append(word_dic[w2][p2])
126             p2+=1
127             if p2==len(word_dic[w2]):
128                 while p1<len(word_dic[w1]):
129                     lst.append(word_dic[w1][p1])
130                     p1+=1
131                 return lst
132         if word_dic[w1][p1]==word_dic[w2][p2]:
133             lst.append(word_dic[w1][p1])
134             p1+=1
135             p2+=1
136             if p1==len(word_dic[w1]):
137                 while p2<len(word_dic[w2]):
138                     lst.append(word_dic[w2][p2])
139                     p2+=1
140                 return lst
141             if p2==len(word_dic[w2]):
142                 while p1<len(word_dic[w1]):
143                     lst.append(word_dic[w1][p1])
144                     p1+=1
145                 return lst
146
147 def mergeNot(w1):
148     lst=[]
149     if w1 not in word_bag:
150         for i in range(1,len(doc_word)):
151             lst.append(i)
152         return lst
153     newlst=[]
154     for i in range(1,len(word_dic[w1])):
155         newlst.append(word_dic[w1][i])
156     for i in range(1,len(doc_word)+1):
157         if i not in word_dic[w1]:
158             lst.append(i)
159     return lst

```

Examples:

```
In [21]: word_dic['harry']
```

```
Out[21]:
```

```
[23,  
 6926,  
 9607,  
12101,  
12175,  
12441,  
12663,  
12700,  
12717,  
12740,  
12788,  
12794,  
12940,  
13602,  
16887,  
16963,  
17034,  
17172,  
18010,  
19124,  
19220,  
19450,  
19511,  
21469]
```

```
In [22]: word_dic['potter']
```

```
Out[22]: [4, 16887, 16963, 17172, 21469]
```

```
In [23]: merge('potter and harry')
```

```
Out[23]: [16887, 16963, 17172, 21469]
```

```
In [24]: merge('potter or harry')
```

```
Out[24]:
```

```
[6926,  
 9607,  
12101,  
12175,  
12441,  
12663,  
12700,  
12717,  
12740,  
12788,  
12794,  
12940,  
13602,  
16887,  
16963,  
17034,  
17172,  
18010,  
19124,  
19220,  
19450,  
19511,  
21469]
```

结论分析与体会：

Using inverted index can help to reduce the space used to store the words, their frequency and the document in which they appear for the matrix is sparse.

Boolean retrieval model can analyze the combination of several logical Operators and return the documents that fit the query.