Report

This is a script which is used to retrieve queries in documents, and find the similarities between them.

To run this script, type command: python retrieval.py (options) documents.txt queries.txt. The options contain:

```
    -s stop_list [whether to use stoplist to filter document]
    -I filename [write results in this file]
    -n number [retrieving for a specific query, number is docid of query]
    -I [whether to use NLTK library to process document]
    -A [retrieving for the full query set]
```

To get the results of full query set, type command:

```
python retrieval.py -s stop_list.txt -l -l results.txt -A documents.txt queries.txt
```

documents.txt and queries.txt will be read by ReadDocuments() in read_documents.py. Then the words are filtered by stoplist and NLTK stemmer. Finally, the frequency of each word will be calculated and stored in a new dictionary with word (word is key, frequency is value).

```
40 #Filter the document
41 def filter document(filename):
       documents = ReadDocuments(filename)
43
       wordRE = re.compile(r'[A-Za-z]+')
       stemmer = PorterStemmer()
44
45
       filtereddoc = []
46
       for doc in documents:
47
           doclist = []
           for line in doc.lines:
48
               line.replace("\n","")
49
               for word in wordRE.findall(line.lower()):
50
51
                   if word not in stops:
                       if '-I' in opts:
52
53
                           word = stemmer.stem(word)
                       doclist.append(word)
54
55
           c = Counter(doclist)
           filtereddoc.append(c)
56
57
       return filtereddoc
52
```

After filtering documents.txt and queries.txt, calculating the tf-idf for each word.

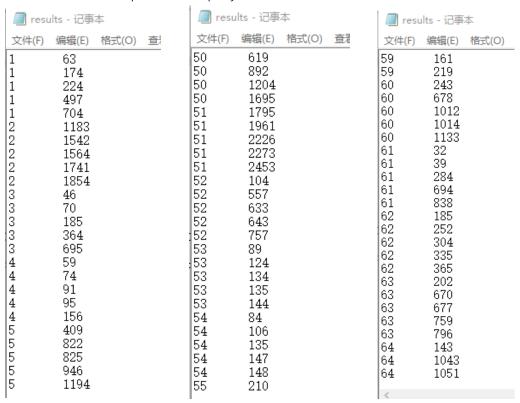
```
60 #calculate idf
61 def tf(word, count):
      return count[word] / sum(count.values())
63
64 def n_containing(word, count_list):
65
      return sum(1 for count in count list if word in count)
67 def idf(word, count list):
      return math.log(len(count list) / (1 + n containing(word, count list)))
70 def tfidf(word, count, count_list):
71
      return tf(word, count) * idf(word, count_list)
73 #filter documents and calculate tf-idf for each document
74 countlist = filter_document(filename)
75 doctfidf = []
76 for i, count in enumerate(countlist):
      scores = {word: tfidf(word, count, countlist) for word in count}
77
78
      doctfidf.append(scores)
79
80 #filter queries and calculate tf-idf for each query
81 querylist = filter_document(queryname)
82 quetfidf = []
83 for i, count in enumerate(querylist):
      scores = {word: tfidf(word, count, querylist) for word in count}
85
      quetfidf.append(scores)
```

Loop the document in queries and documents. For each pair of them, find the same word and calculate cos similarity.

For each document in queries, sort the similarity with document in documents. Then print 5 pairs with the highest similarity into results.txt

```
87 #retrieving all queries
 88 if '-A' in opts:
       for i,query in enumerate(quetfidf):
 89
 90
           sims = \{\}
 91
           for j,doc in enumerate(doctfidf):
 92
                si={}
 93
               for word in quetfidf[i]:
 94
                    if word in doctfidf[j]:
 95
                        si[word] = 1
                if len(si) == 0:
 96
 97
                    sims[j] = 0
 98
                    continue
99
                pSum = sum([quetfidf[i][word]*doctfidf[j][word] for word in si])
100
                sum1Sq = sum([pow(quetfidf[i][word],2) for word in si])
101
                sum2Sq = sum([pow(doctfidf[j][word],2) for word in si])
102
                den = math.sqrt(sum1Sq*sum2Sq)
103
               if den == 0:
104
                    sims[j] = 0
105
                    continue
106
                r = float(pSum)/den
107
                sims[j] = r
108
           sorted_sims = sorted(sims.items(),key=lambda x: x[1],reverse=True)
109
           for k,v in sorted_sims[:5]:
110
                   f.write(str(i+1)+'\t')
111
                    f.write(str(k+1)+'\n')
```

In results.txt, the first column is docid of queries.txt. The second column is docid of documents.txt. Each docid of queries.txt will correspond to 5 docids of documents.txt. They are the first 5 similar pairs in this query.



If type: python retrieval.py -s stop_list.txt -I -I result3.txt -n 3 documents.txt queries.txt It means only retrieve docid 3 in queries.txt. It will compare with each docid in documents.txt. Then sort all the pairs that similarity not equal to 0. The result3.txt:

🧻 result3 - 记事本			
文件(F)	编辑(E)	格式(O)	查看
333333333333333333333333333333333333333	46 70 1163 2189 1173 185 2246 3124 1365 1380 364 1392 2424 1402 1431 3187 2523 2537 2957 1535 1552 1564 1572 2598 2482		