

# Web Lab1 邮件搜索引擎

PB17111585 张永停

## 一、算法

### I、数据处理

本部分代码见 `data_process.py`

- 将每个邮件映射到一个INT数字，映射关系存储为字典，`{Int: EmailAddress}`，存储文件为 `dict.npy`
  - 代码见 `convert_to_dict` 函数

```
1 def convert_to_dict(filelist):
2     i = 0
3     dict = {}
4     for file in filelist:
5         dict[str(i)] = file
6         i = i + 1
7     np.save(r"../dataset/dict.npy", dict)
```

- 分词、词根化、去停用词
  - 代码见 `tokenization` 函数

```
1 def tokenization(dict):
2     dict = np.load(dict, allow_pickle=True).item()
3     snowball_stemmer = SnowballStemmer("english")
4     stopwords = set(stopwords.words('english'))
5     for w in ['!', ',', '.', '?', '-s', '-ly', '</s>', 's']:
6         stopwords.add(w)
7     for key, value in dict.items():
8         #if int(key) < 24502:
9         #    continue
10        dict_fre = {}
11        #print(value)
12        with open(value, 'r', errors='ignore') as f:
13            data = f.read()
14            words = word_tokenize(data)
15            #words_filter = []
16            for word in words:
17                if word not in stopwords:
18                    if snowball_stemmer.stem(word) not in dict_fre:
19                        dict_fre[snowball_stemmer.stem(word)] = 1
20                    else:
21                        dict_fre[snowball_stemmer.stem(word)] += 1
22            path_index = target_path + key
23            np.save(path_index, dict_fre)
24            #t = np.load(path_index + '.npy',
25            allow_pickle=True).item()
26            #print(t)
```

- 此部分需要注意部分邮件还有非ASCII符号，因此读入时设置为 `open(value, 'r', errors='ignore')`
- 采用 `nltk` 进行分词
- 使用 `nltk` 的英文包作为停用词，同时将 `'!', ',', '.?', '-s', '-ly', '</s>', 's'` 加入停用词
- 使用 `nltk` 的 `SnowballStemmer("english")` 来进行词干提取
- 为后续方便，将每个文件的分词结果按形式 `{word, frequency}` 存储在 `\dataset\index` 下，按第一步中的Int命名
- 检索出现频率最高的1000个词
  - 代码见 `get_top1k_words`

```

1  def get_top1k_words():
2      words = {}
3      for file in os.listdir(target_path):
4          word = np.load(target_path + file,
allow_pickle=True).item()
5          for w, f in word.items():
6              if w not in words:
7                  words[w] = f
8              else:
9                  words[w] += f
10         words_order = sorted(words.items(), key=lambda x: x[1],
reverse=True)
11         np.save(top1k_path, words_order[:1000])
12         #print(words_order)

```

- 对上一步分词的存储结果遍历检索，统计词频
- 将前1000的词按照 `{word, frequency}` 存储为 `output\top1k.npy`
- 为后续检索方便，同时将 `{word, index}` 形式字典存储为 `\output\word_index_top1k.npy`
- 建立倒排表
  - 代码见 `get_invered_table`

```

1  def get_invered_table():
2      words_table = []
3      for i in range(1000):
4          words_table.append([])
5          words_fre = np.zeros(1000)
6          top1k = get_dict_top1k()
7          index_word_lastone = [None] * 1000
8          for i in range(517401):
9              #for file in os.listdir(target_path):
10                 top1k_words = {}
11                 word = np.load(target_path + str(i) + '.npy',
allow_pickle=True).item()
12                 for w, f in word.items():

```

```

13         if w in top1k:
14             index_word = list((top1k.keys())).index(w)
15             if words_table[index_word] == []:
16                 words_table[index_word].append(i)
17                 index_word_lastone[index_word] = i
18             else:
19                 lastone = index_word_lastone[index_word]
20                 words_table[index_word].append(i - lastone)
21                 index_word_lastone[index_word] = i
22             top1k_words[w] = f
23             words_fre[index_word] += int(f)
24         #if(i > 5):
25         #    break
26         #np.save(filter_path + str(i), top1k_words)
27
28     #words_fre_path =
29     r'D:\study\USTC\2020.fall\web\lab1\output\freq.npy'
30     #np.save(words_fre_path, words_fre)
31     words_table_path =
32     r'D:\study\USTC\2020.fall\web\lab1\output\invertedtable.npy'
33     np.save(words_table_path, words_table)

```

- 对第二步分词存储的文件遍历，如果文件中的词出现在top1k，则将该文件的INT值存入相应索引的列表中。建立words\_table，形式为[[[]]]，存储于\output\invertedtable.npy
- 这里使用临时数组index\_word\_lastone来记录每个索引上次存入的文件编号，从而在存储时，使用文件编号的间隔代替原编号，节省存储空间。
- 在该步，同时统计每个词出现的次数，并存放入\output\freq.npy

## II、布尔检索

本部分代码见 bool\_search.py

- 处理输入：
  - 对输入的查询进行分词，词干提取，见deal\_search函数

```

1 def deal_search(result_path):
2     global word_index
3     word_index = np.load(word_index_path, allow_pickle=True).item()
4     origin = input("Please input query(Input only E to exit):")
5     while(origin != 'E'):
6         query = origin.split(' ')
7         start = time.time()
8         split_query = []
9         snowball_stemmer = SnowballStemmer("english")
10        for i, q in enumerate(query):
11            # 这里默认不会查询带着 '(' ')' 的词
12            if (len(q) > 1 and q[0] == '('):
13                split_query.append('(')
14                split_query.append(snowball_stemmer.stem(q[1:]))
15            elif (len(q) > 1 and q[-1] == ')'):
16                split_query.append(snowball_stemmer.stem(q[:-1]))
17                split_query.append(')')
18            else:
19                split_query.append(snowball_stemmer.stem(q))

```

```

20
21     vec = bool_search(split_query)
22     result = []
23     dict = np.load(dict_path, allow_pickle=True).item()
24     # print(dict)
25     for i in range(num_file):
26         if vec[i]:
27             result.append(dict[str(i)][48:])
28     end = time.time()
29     print("Finish! See result at %s" % result_path)
30     print("Time:", end - start)
31     with open(result_path, 'a+') as f:
32         f.write("For query: " + origin + "\n")
33         f.write("Result is \n")
34         for r in result:
35             f.write(r)
36             f.write("\n")
37
38     f.write("=====\n")
39     origin = input("Please input query(Input only E to exit):")

```

- 计算

- 本部分采用数组来模拟栈操作，以正确计算 AND 和 NOT 和 OR 以及 ( )，思想借鉴于简单计算器的实现，优先级 NOT > AND > OR，当遇到 )，弹出直到遇到 (
- 采用两个栈: 符号栈与内容栈

上两个部分代码见 bool\_search

```

1  def bool_search(query):
2      global head_oper
3      global head_word
4      for word in query:
5          if(word == '('):
6              head_oper += 1
7              stack_oper[head_oper] = '('
8
9          elif(word == ')'):
10             while(stack_oper[head_oper] != '('):
11                 do_bool(stack_oper[head_oper])
12             head_oper -= 1
13         elif(word == 'not'):
14             head_oper += 1
15             stack_oper[head_oper] = word
16
17         elif(word == 'and'):
18             if head_oper == -1:
19                 head_oper += 1
20                 stack_oper[head_oper] = word
21             elif stack_oper[head_oper] == 'or' or
stack_oper[head_oper] == '(':
22                 head_oper += 1
23                 stack_oper[head_oper] = word
24             elif stack_oper[head_oper] == 'and' or
stack_oper[head_oper] == 'not' :

```

```

25         while(head_oper >= 0 and ([head_oper] == 'and' or
stack_oper[head_oper] == 'not')):
26             do_bool(stack_oper[head_oper])
27             head_oper += 1
28             stack_oper[head_oper] = word
29         elif (word == 'or'):
30             if head_oper == -1 or stack_oper[head_oper] == '(':
31                 head_oper += 1
32                 stack_oper[head_oper] = word
33             else:
34                 while (head_oper >= 0 and ([head_oper] == 'and' or
stack_oper[head_oper] == 'not' or stack_oper[head_oper] == 'or')):
35                     do_bool(stack_oper[head_oper])
36                     head_oper += 1
37                     stack_oper[head_oper] = word
38             else:
39                 head_word += 1
40                 stack_word[head_word] = word_index[word]
41
42     while(head_oper != -1):
43         do_bool(stack_oper[head_oper])
44
45     return stack_word[head_word]

```

- 由于之前在倒排表中存放的是文件间隔，故这里需要恢复文件索引，在恢复文件索引的同时，将倒排表每个索引对应文件变成bool向量

```

1 def get_bool_vec(word):
2     result = np.zeros(num_file).astype(bool)
3     contain_file = inverted_table[word]
4     last_index = 0
5     for f in contain_file:
6         real_index = last_index + f
7         result[real_index] = True
8         last_index += f
9     return result

```

- 对于布尔运算，直接对bool向量进行逻辑运算。代码见 `do_bool`、`get_and` 等

```

1 def do_bool(oper):
2     global head_word
3     global head_oper
4     global result
5     word2 = stack_word[head_word]
6
7     if oper == 'not':
8         stack_word[head_word] = get_not(word2)
9
10    if oper == 'and':
11        word1 = stack_word[head_word - 1]
12        head_word -= 1
13        stack_word[head_word] = get_and(word1, word2)
14    if oper == 'or':
15        word1 = stack_word[head_word - 1]
16        head_word -= 1

```

```

17         stack_word[head_word] = get_or(word1, word2)
18
19     head_oper -= 1
20     return

```

```

1 def get_and(word1, word2):
2     vec1 = word1
3     vec2 = word2
4     if not (type(word1) == type(np.array([0]))):
5         vec1 = get_bool_vec(word1)
6     if not (type(word2) == type(np.array([0]))):
7         vec2 = get_bool_vec(word2)
8     return vec1 & vec2

```

- 结果存储
  - 检索结果存储在 `\output\bool_search_result.txt` 中

### III、语义检索

- 输入处理，本部分输入处理较为简单，只需要对输入的查询进行分词，去停用词，词根提取

```

1 def deal_search(result_path):
2     global tfidf_table
3     tfidf_table = np.load(tfidf_path)
4     dict = np.load(dict_path, allow_pickle=True).item()
5     global word_index
6     global word_idf
7     word_index = np.load(word_index_path, allow_pickle=True).item()
8     word_idf = np.load(top1k_idf_path)
9     global tfidf_norm
10    tfidf_norm = np.load(tfidf_norm_path)
11    snowball_stemmer = SnowballStemmer("english")
12    stopwords = set(stopwords.words('english'))
13    for w in ['!', ',', '.', '?', '-s', '-ly', '</s>', 's']:
14        stopwords.add(w)
15
16    origin = input("Please input query(Input only E to exit):")
17    while(origin != 'E'):
18        start = time.time()
19        query = {}
20        words = word_tokenize(origin)
21        # words_filter = []
22        num_query_words = 0
23        for word in words:
24            if word not in stopwords:
25                num_query_words += 1
26                if snowball_stemmer.stem(word) not in query:
27                    query[snowball_stemmer.stem(word)] = 1
28                else:
29                    query[snowball_stemmer.stem(word)] += 1
30        #print(query)
31        top10_index = semantic_search(query, num_query_words)
32        end = time.time()
33        print("Finish! See result at %s" % result_path)

```

```

34         print("Time:", end - start)
35         with open(result_path, 'a+') as f:
36             f.write("For query: " + origin + "\n")
37             f.write("Result is \n")
38             for index in top10_index:
39                 f.write(dict[str(index)][48:])
40                 f.write("\n")
41
42         f.write("=====
=====\\n")
42         origin = input("Please input query(Input only E to exit):")

```

- tf-idf计算

- 由于矩阵很大，故提前计算好矩阵的tf-idf并存放于硬盘中

```

1  def cac_tfidf():
2      top1k = np.load(top1k_path, allow_pickle=True)
3
4      words_fre = np.zeros(num_index)
5      inverted_table = np.load(inverted_path, allow_pickle=True)
6      for i in range(num_index):
7          words_fre[i] = len(inverted_table[i])
8
9      dict_top1k = {}
10     for i in range(num_index):
11         dict_top1k[top1k[i][0]] = i
12
13
14     for i in range(num_file):
15         word = np.load(filter_path + str(i) + '.npy',
allow_pickle=True).item()
16         for w, f in word.items():
17             if f == 0:
18                 tf = 0
19             else:
20                 tf = 1 + np.log10(f)
21                 index = dict_top1k[w]
22                 idf = np.log10(num_file / words_fre[index])
23                 tfidf_table[index, i] = tf * idf
24     np.save(tfidf_path, tfidf_table)

```

- 同理，对每个关键词计算idf，并存入内存

```

1  def cac_all_idf():
2      top1k = np.load(top1k_path, allow_pickle=True)
3      words_fre = np.zeros(num_index)
4      inverted_table = np.load(inverted_path, allow_pickle=True)
5      for i in range(num_index):
6          words_fre[i] = len(inverted_table[i])
7      words_idf = np.zeros(num_index)
8      for i in range(num_index):
9          words_idf[i] = np.log10(num_file / words_fre[i])
10     np.save(top1k_idf_path, words_idf)

```

- 同时，为了节省时间，提前计算好tf-idf每个向量的模长并存储

```

1 def cac_tfidf_norm():
2     tfidf_table = np.load(tfidf_path)
3     result = np.linalg.norm(tfidf_table, axis=0)
4     np.save(tfidf_norm_path, result)

```

- 检索
  - 使用矩阵计算并返回top10结果

```

1 def semantic_search(query, num_query_words):
2
3     word_vec = np.zeros(num_index)
4
5     #print(word_idf)
6     for w, f in query.items():
7         if w not in word_index:
8             continue
9         index = word_index[w]
10        word_vec[index] = 1 + np.log10(f / num_query_words)
11        word_vec[index] *= word_idf[index]
12    #print(word_vec)
13
14    word_vec = word_vec.reshape((1, num_index))
15    #word_vec = np.array([1,2]).reshape((1,2))
16    #tfidf_table = np.array([[1,2],[3,4]])
17    cos = word_vec.dot(tfidf_table) / (np.linalg.norm(word_vec) *
tfidf_norm)
18    cos = np.squeeze(cos)
19    #print(cos)
20    #print(cos.argsort()[-10:][::-1])
21    top10_index = cos.argsort()[-10:][::-1]
22    print(cos[top10_index])
23    return top10_index

```

## IV、word2vec

- 训练
  - 本部分使用 `gensim` 模块来训练模型
  - 考虑到word2vec的特性，仅进行分词，并不去停用词和提取词干
  - 为了节省内存，使用字节流的形式传入模型

```

1 class MySentences(object):
2     def __init__(self, dict):
3         dict = np.load(dict, allow_pickle=True).item()
4         self.dict = dict
5
6     def __iter__(self):
7         for key, value in self.dict.items():
8             with open(value, 'r', errors='ignore') as f:
9                 data = f.read()
10                words = word_tokenize(data)
11                #print(words)
12                #break
13            yield words
14

```



```

15 def train():
16     sentences = MySentences(dict_path)
17     #sentences.__iter__()
18     model = gensim.models.word2vec(sentences)
19     model.save(model_path)

```

- 计算每个文件的向量

- 将每个文件进行分词，然后传入模型得到每个词的向量，相加取平均的结果作为该文件的向量

```

1  def sentence_vector(model, words):
2      l = 0
3      v = np.zeros(100)
4      for word in words:
5          if word in model:
6              v += model[word]
7              l += 1
8      if l == 0:
9          return np.zeros(100)
10     v /= l
11     return v
12
13 def cac_file_vec():
14     model = gensim.models.word2vec.load(model_path)
15     file_vec = [[] for i in range(517401)]
16     #print(model['offer'].shape)
17     dict = np.load(dict_path, allow_pickle=True).item()
18     for key, value in dict.items():
19         with open(value, 'r', errors='ignore') as f:
20             data = f.read()
21             words = word_tokenize(data)
22
23             file_vec[int(key)] = sentence_vector(model, words)
24
25     np.save(file_vec_path, file_vec)

```

- 同时，同语义检索部分，预先将所有向量的模计算并存入磁盘

```

1  def cac_vec_file_norm():
2      file_vec = np.array(np.load(file_vec_path, allow_pickle=True)
3                           [:517401])
4      temp = []
5      for arr in file_vec:
6          temp.append(arr)
7      file_vec = np.array(temp)
8      #print(file_vec)
9      file_vec = file_vec.transpose()
10     np.save(file_vec_numpy_path, file_vec)
11     #print(file_vec.shape)
12     result = np.linalg.norm(file_vec, axis=0)
13     #print(result)
14     np.save(file_vec_norm_path, result)

```

- 检索

- 检索过程与语义检索极其相似，只是这里使用word2vec的模型得到查询向量，然后和所有文件计算余弦相似度，并取top10

```

1 def word2vec_search(vec):
2     vec = vec.reshape((1, 100))
3     cos = vec.dot(file_vec) / (np.linalg.norm(vec) * file_vec_norm)
4     cos = np.squeeze(cos)
5     top10_index = cos.argsort()[-10:][::-1]
6     #print(cos[top10_index])
7     return top10_index

```

- 输入处理

- 输入与语义检索部分几乎完全相同，只是不去除停用词，提取词根

```

1 def deal_search(result_path):
2     model = gensim.models.word2vec.load(model_path)
3     global file_vec
4     file_vec = np.load(file_vec_numpy_path)
5     global file_vec_norm
6     file_vec_norm = np.load(file_vec_norm_path)
7     dict = np.load(dict_path, allow_pickle=True).item()
8     origin = input("Please input query(Input only E to exit):")
9     while (origin != 'E'):
10         start = time.time()
11         words = word_tokenize(origin)
12         query = sentence_vector(model, words)
13
14         top10_index = word2vec_search(query)
15         end = time.time()
16         print("Finish! See result at %s" % result_path)
17         print("Time:", end - start)
18         with open(result_path, 'a+') as f:
19             f.write("For query: " + origin + "\n")
20             f.write("Result is \n")
21             for index in top10_index:
22                 f.write(dict[str(index)][48:])
23                 f.write("\n")
24
25         f.write("=====\n")
26         origin = input("Please input query(Input only E to exit):")

```

## 二、优化

### I、压缩倒排表

采用间距来代替文档ID，这部分详见 一、II、的倒排表部分

### II、word2vec

该部分详见 一、IV

结果见下

## 三、实验结果

选取 power、price、meeting、issues、offer

## I、布尔查询

- power and price

- Time: 0.26488184928894043

- 由于结果很多，结果存放于 output/bool\_search\_result 中，助教可以查看~~~

- | For query: power and price

- Result is

- maildir\allen-p\all\_documents\10

- maildir\allen-p\all\_documents\19

- maildir\allen-p\all\_documents\241

- maildir\allen-p\all\_documents\261

- maildir\allen-p\all\_documents\359

- maildir\allen-p\all\_documents\379

- maildir\allen-p\all\_documents\380

- not meeting

- Time: 0.4627113342285156

- For query: not meeting

- Result is

- maildir\allen-p\all\_documents\1

- maildir\allen-p\all\_documents\102

- maildir\allen-p\all\_documents\103

- maildir\allen-p\all\_documents\105

- maildir\allen-p\all\_documents\107

- maildir\allen-p\all\_documents\108

- maildir\allen-p\all\_documents\109

- not issues

- Time: 0.46822571754455566

- For query: not issues

- Result is

- maildir\allen-p\all\_documents\10

- maildir\allen-p\all\_documents\102

- maildir\allen-p\all\_documents\103

- maildir\allen-p\all\_documents\104

- maildir\allen-p\all\_documents\105

- maildir\allen-p\all\_documents\106

- maildir\allen-p\all\_documents\107

- maildir\allen-p\all\_documents\108

- issues or meeting

- Time: 0.3431692123413086

- For query: issues or meeting

- Result is

- maildir\allen-p\all\_documents\1

- maildir\allen-p\all\_documents\10

- maildir\allen-p\all\_documents\100

- maildir\allen-p\all\_documents\101

- maildir\allen-p\all\_documents\104

- maildir\allen-p\all\_documents\106

- maildir\allen-p\all\_documents\114

- maildir\allen-p\all\_documents\115

- (power and price) and not issues or offer and meeting

- Time: 0.32477688789367676

```

For query: (power and price) and not issues or offer and meeting
Result is
maildir\allen-p\all_documents\10
maildir\allen-p\all_documents\19
o maildir\allen-p\all_documents\241
maildir\allen-p\all_documents\261
maildir\allen-p\all_documents\359
maildir\allen-p\all_documents\379
maildir\allen-p\all_documents\380

```

## II、语义检索

- power price

- o Time: 0.2953352928161621

```

For query: power price
Result is
maildir\weldon-c\stagecoach\13
maildir\weldon-c\all_documents\194
o maildir\weldon-c\discussion_threads\103
maildir\lavorato-j\all_documents\390
maildir\white-s\deleted_items\930
maildir\white-s\deleted_items\929

```

- meeting price

- o Time: 0.2191011905670166

```

For query: meeting price
Result is
maildir\shively-h\calendar\13
maildir\mckay-j\tasks\4
o maildir\rogers-b\inbox\1
maildir\rogers-b\all_documents\1766
maildir\campbell-l\inbox\230
maildir\rogers-b\sent\714
maildir\campbell-l\all_documents\1795
maildir\campbell-l\all_documents\1258

```

- issues meeting offer

- o Time: 0.21761775016784668

```

For query: issues meeting offer
Result is
maildir\shackleton-s\meetings\30
maildir\shackleton-s\all_documents\11376
o maildir\mckay-b\all_documents\65
maildir\mckay-b\sent\40
maildir\hyvl-d\all_documents\1015
maildir\campbell-l\all_documents\1083
maildir\hyvl-d\gas\hpl_customers\57

```

- power offer issues

- o Time: 0.22281527519226074

- 
- **For query:** power offer issues  
Result is  
maildir\ring-r\eesirenewableenergy\22  
maildir\presto-k\deleted\_items\636
    - maildir\lavorato-j\sent\_items\102  
maildir\sager-e\calendar\32  
maildir\mckay-b\all\_documents\65  
maildir\kaminski-v\deleted\_items\1169  
maildir\dorland-c\sent\_items\573
  - power price meeting issues offer
    - Time: 0.24132943153381348  
**For query:** power price meeting issues offer  
Result is  
maildir\taylor-m\all\_documents\3936  
maildir\quigley-d\deleted\_items\208
      - maildir\rogers-b\inbox\1  
maildir\weldon-c\stagecoach\13  
maildir\weldon-c\misc\_gas\_issues\3  
maildir\weldon-c\misc\_gas\_issues\1  
maildir\kean-s\heat\_wave\155

### III、word2vec

- power price
    - Time: 0.09765958786010742  
**For query:** power price  
Result is  
maildir\kean-s\discussion\_threads\1789  
maildir\kean-s\all\_documents\2141  
maildir\kean-s\calendar\untitled\2142  
maildir\kean-s\archiving\untitled\186
      - maildir\dasovich-j\notes\_inbox\3712  
maildir\dasovich-j\all\_documents\11968  
maildir\dasovich-j\notes\_inbox\2325  
maildir\dasovich-j\all\_documents\4380  
maildir\dasovich-j\all\_documents\4336  
maildir\dasovich-j\notes\_inbox\2284
  - meeting price
    - Time: 0.07728171348571777
-

For query: meeting price

Result is

maildir\crandell-s\inbox\rto\_west\12  
maildir\campbell-l\discussion\_threads\1615  
maildir\campbell-l\notes\_inbox\567  
○ maildir\campbell-l\all\_documents\1783  
maildir\jones-t\sent\5309  
maildir\jones-t\all\_documents\8974  
maildir\hain-m\discussion\_threads\633  
maildir\hain-m\notes\_inbox\538  
maildir\hain-m\all\_documents\647  
maildir\steffes-j\ees\_mgmt\6

- issues meeting offer

- Time: 0.07729029655456543

For query: issues meeting offer

Result is

maildir\kean-s\all\_documents\677  
maildir\kean-s\discussion\_threads\614  
maildir\kean-s\archiving\untitled\1651  
○ maildir\kean-s\calendar\untitled\677  
maildir\kean-s\wto\6  
maildir\hain-m\all\_documents\99  
maildir\kean-s\all\_documents\513  
maildir\hain-m\discussion\_threads\1141  
maildir\kean-s\sent\53  
maildir\kean-s\archiving\untitled\1814

- power offer issues

- Time: 0.07552051544189453

For query: power offer issues

Result is

maildir\dasovich-j\all\_documents\7745  
maildir\dasovich-j\notes\_inbox\6124  
maildir\kean-s\all\_documents\1743  
○ maildir\kean-s\discussion\_threads\1488  
maildir\kean-s\archiving\untitled\584  
maildir\kean-s\calendar\untitled\1744  
maildir\fossum-d\all\_documents\616  
maildir\kean-s\heat\_wave\29  
maildir\fossum-d\discussion\_threads\369  
maildir\fossum-d\notes\_inbox\334

- power price meeting issues offer

- Time: 0.07729864120483398

For query: power price meeting issues offer

Result is

maildir\kean-s\discussion\_threads\1488  
maildir\kean-s\all\_documents\1743  
maildir\kean-s\heat\_wave\29  
○ maildir\kean-s\archiving\untitled\584  
maildir\kean-s\calendar\untitled\1744  
maildir\fossum-d\all\_documents\616  
maildir\fossum-d\discussion\_threads\369  
maildir\fossum-d\notes\_inbox\334  
maildir\scott-s\discussion\_threads\423  
maildir\scott-s\all\_documents\524

可以看到word2vec速度明显快于语义检索

## 四、运行环境与方法

- python 3.7
  - nltk、numpy、gensim
- bool检索运行方法
  - 在 `bool_search.py` 最下更改希望结果存储的path, 然后 `python3 bool_search.py`
  - 根据提示输入bool查询

```
Please input query(Input only E to exit):(power and price) and not
issues or offer and meeting
Finish! See result at D:\study\USTC\2020
.fall\web\lab1\output\bool_search_result.txt
Time: 0.32477688789367676
Please input query(Input only E to exit):E

Process finished with exit code 0
```

- 语义检索运行方法
  - 在 `semantic_search.py` 最下更改希望结果存储的path, 然后 `python3 semantic_search.py`
  - 根据提示输入查询

```
Please input query(Input only E to exit):offer
[0.46981613 0.46431106 0.45986456 0.43528531 0.3991407 0.38775676
0.3852518 0.38375876 0.3829379 0.38075948]
Finish! See result at D:\study\USTC\2020
.fall\web\lab1\output\semantic_search_result.txt
Time: 0.20234274864196777
Please input query(Input only E to exit):E
```

- word2vec运行方法
  - 在 `word2vec.py` 最下更改希望结果存储的path, 然后 `python3 word2vec.py`
  - 根据提示输入查询

```
Please input query(Input only E to exit):power price meeting issues offer
Finish! See result at D:\study\USTC\2020.fall\web\lab1\output\word2vec_search_result.txt
Time: 0.07729864120483398
Please input query(Input only E to exit):E

Process finished with exit code 0
```