

Web信息处理与应用 实验报告

lab1

团队成员

- 何春望 PB17000075
- 吴健宗 PB17000082

实验内容

本实验要求以给定的邮件数据集为基础，实现一个邮件搜索引擎。对于给定的查询，能够以**精确查询**和**模糊语义匹配**的方式返回最相关的一系列邮件文档。

算法描述

- 精确查询 (bool 查询)：对于给定的 bool 查询，如 (term1 AND term2 OR NOT term3)，根据倒排索引表，返回符合查询的所有文档。
- 模糊语义匹配：对文档集合建立 tf-idf 矩阵，对输入的语义查询 Q，计算 Q 的 tf-idf 向量 v ，并找出 v 与文档集合中的文档相似度最高的文档子集。

优化内容

时间复杂度优化

- 跳表指针
 - bool 查询中，对 AND，OR 的处理，使用了**跳表指针**技术，减少了时间复杂度。

空间复杂度优化

- 减小 tf-idf 矩阵精度
 - 对 tf-idf 矩阵的设置，从 numpy 默认的 float64 精度减小到 float16 精度。因为在文档数量，单词数量有限的情况下，tf-idf 矩阵的可能值实际上是有限的、离散的，不需要很高的精度就可以达到筛选文档相似度最高的文档子集的目标。这样做减少了 tf-idf 矩阵存储的空间复杂度。
 - 对所有文档建立 tf-idf 矩阵时，如果使用 float16 的精度，只需要 900MB-1000MB 的空间，如果使用 float64，则需要四倍左右的空间。

注：实验结果表示，在倒排索引，tf-idf 矩阵，单词索引三个数据库中，tf-idf 矩阵占用空间远大于其他两个，因此减少 tf-idf 矩阵的空间大小是最有必要的。

- 去邮件头处理
 - 针对 enron dataset 的文档格式，我们做了特殊的“去邮件头”处理。具体是：对每一个邮件，先去掉一些对查询帮助不大的内容，如：邮件发送人，接受人，邮件版本 (Mime-Version)，日期，邮件编号等。这样做可以减少内存消耗，并大大加快后期的处理速度，同时减小倒排索引表，tf-idf 矩阵等数据库的大小。
 - 经粗略统计，enron dataset 中，每条邮件大概有三分之二是“无效信息”，因此先把它们去掉是很有帮助的。

注：去邮件头处理在去除邮件编号，邮件版本等“冗余信息”的同时，貌似还去掉了日期，发件人等在实际应用中可能有效的信息。考虑到本实验的查询词是针对邮件正文进行查询，因此在本实验的查询词范围内，这些信息是无效的；如果需要查询这些信息，在去邮件头处理中不对它们做去除处理即可。

- 去停用词处理
 - 本实验中，我们采用了去停用词处理以减少倒排表和 tf-idf 文档的大小。

- 我们同时采用了**两个停用词表**，一个是通用英文停用词表，包含 I, you, do 等常用停用词。另一个是针对 enron dataset 的特殊停用词表，包含 enron, com, html 等在该数据集中频繁出现的无效词语。

查询结果优化

- **词根化处理**
 - 本实验中，我们调用了 nltk 库对邮件中的单词进行词根化。这样做有利于避免因为词性转化而导致搜索不到目标词的问题。如，用户需要搜索 "message"，但某个文档中出现的是 "messages"，即 "message" 的复数形式。如果不做词根化处理，搜索引擎则不会把他们识别为同一个词，导致遗漏的情况。
 - 另外，我们实现了三种不同的词根化方案，可以满足对构建时间和准确度不同的权衡与取舍要求，极大地提高了系统的灵活性。


实验结果

- 构建索引时间：1h 4min
- 查询延迟：
 - bool 检索：<1s
 - 语义检索：~7s
- 倒排索引表占用空间：164MB
- tf-idf 矩阵占用空间：986MB

运行示例

build inv-index and tf-idf matrix

运行 build_index_and_matrix.py，大约一个小时候，得到结果：



```
Windows PowerShell
300/1000 word tfidf processed
400/1000 word tfidf processed
500/1000 word tfidf processed
600/1000 word tfidf processed
700/1000 word tfidf processed
800/1000 word tfidf processed
900/1000 word tfidf processed
1000/1000 word tfidf processed
pleas 232689
power 60804
energi 63858
market 71675
compani 61711
time 125019
thank 191586
messag 129729
mail 86269
price 59771
['pleas' 'power' 'energi' 'market' 'compani' 'time' 'thank' 'messag'
'mail' 'price' 'corp' 'gas' 'forward' 'call' 'inform' 'day' 'trade'
'origin' 'week' 'servic' 'meet' 'busi' 'deal' 'report' 'california'
'develop' 'schedul' 'chang' 'attach' 'imag']
[[0.617 0. 0. ... 0. 0. 0.3472]
[0. 2.354 0. ... 0. 0. 0. ]
[0.9087 1.817 0. ... 0. 0. 0. ]
[1.678 1.753 0.8584 ... 0. 0. 0. ]
[0. 0.9233 0. ... 0. 0. 0. ]]
build time used: 3861.0709948539734 s
PS C:\Users\hechu\git\enron-search-engine\src>
```

以下是文档频率排行前十的结果，可以看到，基本去除了冗余的**停用词**。

```
pleas 232689
power 60804
energi 63858
market 71675
compani 61711
time 125019
thank 191586
messag 129729
mail 86269
price 59771
```

bool search

- 选用查询词: power, company, customers, employees, president
- 查询用例

注: 查询结果中文档的排序是按照文件名字典序排序的

- **power AND company AND customers AND employees AND president**

执行结果:

```
Windows PowerShell
PS C:\Users\hechu\git\enron-search-engine\src> python .\bool_search.py
Please input your bool query (words connected with AND, OR, NOT): power AND company AND customers
AND employees AND president
bool search time: 0.34359192848205566 s
total number of results: 1816
first 100 results: [3656, 3657, 3917, 3933, 4230, 4251, 4386, 4393, 4415, 4442, 4512, 4548, 4593
, 4761, 4767, 4812, 4814, 5086, 5091, 5191, 5192, 5337, 5371, 5388, 5508, 5509, 6117, 6132, 7016,
7687, 7703, 8801, 9229, 9273, 9955, 11173, 12399, 12579, 14254, 14701, 18318, 19338, 21799, 2344
5, 24296, 24514, 25755, 28503, 28511, 28515, 28565, 28939, 29005, 32544, 32827, 32881, 32913, 329
85, 33011, 33270, 33275, 33591, 33691, 33779, 33786, 35825, 36870, 37394, 38939, 40066, 40178, 41
787, 43249, 43465, 43995, 44027, 44579, 44646, 44798, 44908, 46814, 49117, 49122, 49123, 49132, 4
9135, 49355, 49436, 49501, 49506, 49507, 49516, 49532, 49534, 51966, 52212, 52315, 52736, 53498,
54146]
first ten file path:
('..\dataset\maildir\arnold-j\all_documents\605',)
('..\dataset\maildir\arnold-j\all_documents\606',)
('..\dataset\maildir\arnold-j\all_documents\846',)
('..\dataset\maildir\arnold-j\all_documents\860',)
('..\dataset\maildir\arnold-j\deleted_items\179',)
('..\dataset\maildir\arnold-j\deleted_items\198',)
('..\dataset\maildir\arnold-j\deleted_items\323',)
('..\dataset\maildir\arnold-j\deleted_items\33',)
('..\dataset\maildir\arnold-j\deleted_items\35',)
('..\dataset\maildir\arnold-j\deleted_items\374',)
PS C:\Users\hechu\git\enron-search-engine\src>
```

可以看到, 对五个词用AND连起来查询, 一共用了 0.34s, 查询到 1816 个结果。

下面是第一个结果的部分内容(..\dataset\maildir\arnold-j\all_documents/605):

```
605
dataset > maildir > arnold-j > all_documents > 605
1 Message-ID: <19189634.1075857603655.JavaMail.evans@thyme>
2 Date: Mon, 14 May 2001 01:30:00 -0700 (PDT)
3 From: ann.schmidt@enron.com
4 Subject: Enron Mentions - 05/12/01 - 05/13/01
5 Mime-Version: 1.0
6 Content-Type: text/plain; charset=us-ascii
7 Content-Transfer-Encoding: 7bit
8 X-From: Ann M Schmidt
9 X-To:
10 X-cc:
11 X-bcc:
12 X-Folder: \John_Arnold_Jun2001\Notes Folders\All documents
13 X-Origin: Arnold-J
14 X-FileName: Jarnold.nsf
15
16 As Final Exams Begin, Power Is a Big Question
17 The New York Times, 05/13/01
18
19 British Telecom
20 The Times of London, 05/12/01
21
22 Houston needs to think small about future technology
23 Houston Chronicle, 05/13/01
24
25 Panel plots new course for area's future / Education, economics, quality of
26 life top group's list of needed improvements
27 Houston Chronicle, 05/13/01
28
29 MSEB not to pick up 15 pc in DPC after phase II completion
30 Press Trust of India Limited, 05/13/01
31
32 Enron plans to pull out of Gulf gas project: MEED
33 Agence France-Presse, 05/13/01
34
35 SMALL BUSINESS / Pleasure cruisin' / Yacht fleet owner offers customers what
36 amounts to limo service on the lake
37 Houston Chronicle, 05/13/01
38
```

这是一封长邮件，里面涵盖了全部五个搜索词，正因为搜索词要求五个词都在，所以几乎只有这种长邮件才满足要求。

- **power AND company OR customers AND employees OR president**

执行结果：

```
Windows Powershell
PS C:\Users\hechu\git\enron-search-engine\src> python .\bool_search.py
Please input your bool query (words connected with AND, OR, NOT): power AND company OR customers
AND employees OR president
bool search time: 0.34718823432922363 s
total number of results: 37908
first 100 results: [1, 2, 8, 16, 17, 51, 148, 156, 280, 339, 385, 406, 411, 486, 505, 509, 527,
617, 633, 634, 648, 658, 711, 714, 734, 751, 770, 773, 775, 786, 788, 805, 819, 925, 930, 942, 94
4, 951, 1019, 1049, 1050, 1055, 1066, 1096, 1128, 1138, 1185, 1198, 1201, 1214, 1265, 1269, 1318,
1353, 1371, 1424, 1433, 1442, 1448, 1468, 1470, 1471, 1506, 1524, 1616, 1624, 1788, 1854, 1858,
1927, 1931, 1951, 1953, 1973, 2029, 2033, 2042, 2044, 2086, 2122, 2244, 2254, 2255, 2307, 2405, 2
451, 2458, 2459, 2496, 2589, 2597, 2753, 2819, 2823, 2829, 2897, 2907, 2916, 2919, 2937]
first ten file path:
('..\dataset\maildir\allen-p\all_documents\1',)
('..\dataset\maildir\allen-p\all_documents\10',)
('..\dataset\maildir\allen-p\all_documents\105',)
('..\dataset\maildir\allen-p\all_documents\112',)
('..\dataset\maildir\allen-p\all_documents\113',)
('..\dataset\maildir\allen-p\all_documents\146',)
('..\dataset\maildir\allen-p\all_documents\234',)
('..\dataset\maildir\allen-p\all_documents\241',)
('..\dataset\maildir\allen-p\all_documents\355',)
('..\dataset\maildir\allen-p\all_documents\408',)
PS C:\Users\hechu\git\enron-search-engine\src>
```


将其中两个 AND 换成 OR，搜索结果数量增加为 37908 个。

```
dataset > maildir > allen-p > all_documents > 1
51 Corner
52
53 5. HOT REPORT: Oscar Gruss & Son's most recent issue of its Broadband
54 Brief reports the latest developments in the broadband space.
55
56 6. EDITOR'S PICK: Bear Stearns measures the impact of broadband and the
57 Internet on telecom in Latin America.
58
59 7. FREE STOCK SNAPSHOT: The current analysts' consensus rates 3M (MMM), a
60 "buy/hold."
61
62 8. JOIN THE MARKETBUZZ: where top financial industry professionals answer
63 your questions and offer insights every market day from noon 'til 2:00
64 p.m. ET.
65
66 9. TRANSCRIPTS FROM WALL STREET: Ash Rajan, senior vice president and
67 market analyst with Prudential Securities, answers questions about the
68 market.
69
70 ===== Sponsored by =====
71 Profit From AAI's "Cash Rich" Stock Screen - 46% YTD Return
72
```

查看第一个文档，其中含有 president 这个词，结果正确。

- **power AND company AND NOT (customers OR employees OR president)**

执行结果：

```
Windows Powershell
PS C:\Users\hechu\git\enron-search-engine\src> python .\bool_search.py
Please input your bool query (words connected with AND, OR, NOT): power AND company AND NOT (cust
omers OR employees OR president)
bool search time: 0.5329921245574951 s
total number of results: 8172
first 100 results: [156, 486, 775, 819, 944, 1214, 1624, 1931, 2086, 2244, 2254, 2255, 2597, 289
7, 3157, 3289, 3493, 3660, 4232, 4284, 4299, 4302, 4399, 4519, 4594, 4813, 4976, 5189, 5228, 5392
, 5512, 5745, 6213, 6352, 6525, 6632, 6677, 6837, 6852, 7316, 7567, 7921, 8008, 8010, 8017, 8018,
8020, 8062, 8086, 8097, 8125, 8129, 8140, 8172, 8181, 8183, 8324, 8342, 8489, 8559, 8612, 8619,
8627, 8731, 8780, 8821, 9036, 9044, 9051, 9138, 9156, 9264, 9265, 9267, 9315, 9333, 9458, 9488, 9
489, 9490, 9494, 9517, 9520, 9554, 9555, 9556, 9558, 9560, 9561, 9580, 9584, 9585, 9640, 9641, 96
45, 9646, 9648, 9678, 9731, 9735]
first ten file path:
('..\dataset\maildir\allen-p\all_documents\241',)
('..\dataset\maildir\allen-p\all_documents\540',)
('..\dataset\maildir\allen-p\deleted_items\238',)
('..\dataset\maildir\allen-p\deleted_items\344',)
('..\dataset\maildir\allen-p\deleted_items\52',)
('..\dataset\maildir\allen-p\discussion_threads\427',)
('..\dataset\maildir\allen-p\sent\195',)
('..\dataset\maildir\allen-p\sent\473',)
('..\dataset\maildir\allen-p\sent_items\104',)
('..\dataset\maildir\allen-p\sent_items\248',)
PS C:\Users\hechu\git\enron-search-engine\src>
```

含有 power 和 company 的同时不含有后面三个词，查询结果有 8172 个。

```
dataset > maildir > allen-p > all_documents > 241
1 Message-ID: <3536013.1075855670762.JavaMail.evans@thyme>
2 Date: Thu, 22 Jun 2000 00:26:00 -0700 (PDT)
3 From: phillip.allen@enron.com
4 To: scott.carter@chase.com
5 Subject: Re: The New Power Company
6 Mime-Version: 1.0
7 Content-Type: text/plain; charset=us-ascii
8 Content-Transfer-Encoding: 7bit
9 X-From: Phillip K Allen
10 X-To: "Carter, Scott" <Scott.Carter@chase.com> @ ENRON
11 X-cc:
12 X-bcc:
13 X-Folder: \Phillip_Allen_Dec2000\Notes Folders\All documents
14 X-Origin: Allen-P
15 X-FileName: pallen.nsf
16
17 Scott,
18
19 I emailed your question to a friend that works for the new company. I think
20 I know the answer to your questions but I want to get the exact details from
21 him. Basically, they will offer energy online at a fixed price or some
22 price that undercuts the current provider. Then once their sales are large
23 enough they will go to the wholesale market to hedge and lock in a profit.
24 The risk is that they have built in enough margin to give them room to manage
25 the price risk. This is my best guess. I will get back to you with more.
26
27 Phillip
```

第一个文档中，含有 Power, Company, 同时不含有其他三个词。

- **NOT power AND NOT company AND NOT customers AND NOT employees AND NOT president**

执行结果：

```
Windows Powershell
PS C:\Users\hechu\git\enron-search-engine\src> python .\bool_search.py
Please input your bool query (words connected with AND, OR, NOT): NOT power AND NOT company AND N
OT customers AND NOT employees AND NOT president
bool search time: 1.000763177871704 s
total number of results: 373794
first 100 results: [5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 18, 19, 20, 21, 23, 24, 25, 26, 27, 28,
30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 45, 46, 47, 49, 50, 52, 53, 56, 57, 58, 6
1, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85
, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 100, 101, 102, 103, 104, 105, 106, 107, 108
, 109, 110, 111, 112, 113, 114, 115, 116, 117]
first ten file path:
('..\dataset\maildir\allen-p\all_documents\102',)
('..\dataset\maildir\allen-p\all_documents\103',)
('..\dataset\maildir\allen-p\all_documents\104',)
('..\dataset\maildir\allen-p\all_documents\106',)
('..\dataset\maildir\allen-p\all_documents\107',)
('..\dataset\maildir\allen-p\all_documents\108',)
('..\dataset\maildir\allen-p\all_documents\109',)
('..\dataset\maildir\allen-p\all_documents\11',)
('..\dataset\maildir\allen-p\all_documents\110',)
('..\dataset\maildir\allen-p\all_documents\111',)
PS C:\Users\hechu\git\enron-search-engine\src>
```

搜索不含有这五个词的结果，一共有 373794 个。

- **power OR company OR customers OR employees OR president**

执行结果：

```
Windows Powershell
PS C:\Users\hechu\git\enron-search-engine\src> python .\bool_search.py
Please input your bool query (words connected with AND, OR, NOT): power OR company OR customers O
R employees OR president
bool search time: 0.3810863494873047 s
total number of results: 143607
first 100 results: [1, 2, 3, 4, 8, 16, 17, 22, 29, 44, 48, 51, 54, 55, 59, 60, 99, 137, 143, 148
, 149, 156, 165, 178, 196, 198, 244, 246, 248, 264, 265, 270, 275, 278, 279, 280, 281, 282, 283,
284, 294, 295, 296, 297, 299, 306, 308, 310, 315, 316, 321, 322, 329, 339, 340, 352, 354, 362, 37
2, 385, 390, 392, 400, 401, 406, 411, 412, 418, 425, 426, 432, 443, 447, 472, 481, 486, 487, 492,
504, 505, 509, 511, 520, 527, 535, 556, 563, 583, 602, 617, 625, 626, 633, 634, 635, 637, 642, 6
43, 644, 647]
first ten file path:
('..\dataset\maildir\allen-p\all_documents\1',)
('..\dataset\maildir\allen-p\all_documents\10',)
('..\dataset\maildir\allen-p\all_documents\100',)
('..\dataset\maildir\allen-p\all_documents\101',)
('..\dataset\maildir\allen-p\all_documents\105',)
('..\dataset\maildir\allen-p\all_documents\112',)
('..\dataset\maildir\allen-p\all_documents\113',)
('..\dataset\maildir\allen-p\all_documents\118',)
('..\dataset\maildir\allen-p\all_documents\125',)
('..\dataset\maildir\allen-p\all_documents\14',)
PS C:\Users\hechu\git\enron-search-engine\src>
```

至少包含一个搜索词的结果，有 143607 个。

semantic search

- 选用查询词：opportunity, management, address, offer, price
- 查询用例
 - **opportunity**

执行结果：

```
Windows Powershell
PS C:\Users\hechu\git\enron-search-engine\src> python .\semantic_search.py
Please input your semantic query (words seperated with space): opportunity
('..\dataset\maildir\zipper-a\sent_items\63',) cosine similarity: 0.7954264524103832
('..\dataset\maildir\arnold-j\deleted_items\396',) cosine similarity: 0.6296477495107632
('..\dataset\maildir\mccconnell-m\all_documents\210',) cosine similarity: 0.49047256097560976
('..\dataset\maildir\mccconnell-m\sent\80',) cosine similarity: 0.49047256097560976
('..\dataset\maildir\mccconnell-m\sent_mail\79',) cosine similarity: 0.49047256097560976
('..\dataset\maildir\zipper-a\deleted_items\36',) cosine similarity: 0.46243093922651934
('..\dataset\maildir\ybarbo-p\inbox\184',) cosine similarity: 0.3527960526315789
('..\dataset\maildir\cuilla-m\sent_items\115',) cosine similarity: 0.2903880866425993
('..\dataset\maildir\arnold-j\sent_items\578',) cosine similarity: 0.2783304498269896
('..\dataset\maildir\cuilla-m\sent_items\116',) cosine similarity: 0.26790174854288096
semantic search time: 7.331348180770874 s
PS C:\Users\hechu\git\enron-search-engine\src>
```

输出十个匹配度最高的结果。


```
63
dataset > maildir > zipper-a > sent_items > 63
1  Message-ID: <18458479.1075845422375.JavaMail.evans@thyme>
2  Date: Thu, 24 May 2001 08:02:01 -0700 (PDT)
3  From: andy.zipper@enron.com
4  To: greg.piper@enron.com
5  Subject: RE: possible opportunity with Enron
6  Mime-Version: 1.0
7  Content-Type: text/plain; charset=us-ascii
8  Content-Transfer-Encoding: 7bit
9  X-From: Zipper, Andy </O=ENRON/OU=NA/CN=RECIPIENTS/CN=AZIPPER>
10 X-To: Piper, Greg </O=ENRON/OU=NA/CN=RECIPIENTS/CN=Gpiper>
11 X-cc:
12 X-bcc:
13 X-Folder: \Zipper, Andy\Zipper, Andy\Sent Items
14 X-Origin: ZIPPER-A
15 X-FileName: Zipper, Andy.pst
16
17 yadda yadda yadda....
18
```

查看第一个结果，短邮件，含有 opportunity.

o opportunity, management

执行结果：

```
Windows Powershell
PS C:\Users\hechu\git\enron-search-engine\src> python .\semantic_search.py
Please input your semantic query (words seperated with space): opportunity management
('..../dataset/maildir\\lavorato-j\\all_documents\\698',) cosine similarity: 1.0770101925254814
('..../dataset/maildir\\lavorato-j\\sent_items\\48',) cosine similarity: 1.0770101925254814
('..../dataset/maildir\\lavorato-j\\sent\\574',) cosine similarity: 1.0770101925254814
('..../dataset/maildir\\lavorato-j\\deleted_items\\190',) cosine similarity: 1.0770101925254814
('..../dataset/maildir\\zipper-a\\sent_items\\63',) cosine similarity: 0.7954264524103832
('..../dataset/maildir\\arnold-j\\deleted_items\\396',) cosine similarity: 0.6296477495107632
('..../dataset/maildir\\sanders-r\\deleted_items\\477',) cosine similarity: 0.503707627118644
('..../dataset/maildir\\mcconnell-m\\sent\\80',) cosine similarity: 0.49047256097560976
('..../dataset/maildir\\mcconnell-m\\all_documents\\210',) cosine similarity: 0.49047256097560976
('..../dataset/maildir\\mcconnell-m\\_sent_mail\\79',) cosine similarity: 0.49047256097560976
semantic search time: 7.339157581329346 s
PS C:\Users\hechu\git\enron-search-engine\src> |
```

分别展现第一和第八的结果：


```
dataset > maildir > lavorato-j > all_documents > 698
1 Message-ID: <16305534.1075845545222.JavaMail.evans@thyme>
2 Date: Wed, 27 Dec 2000 23:42:00 -0800 (PST)
3 From: john.lavorato@enron.com
4 To: rick.buy@enron.com, ted.murphy@enron.com
5 Subject:
6 Mime-Version: 1.0
7 Content-Type: text/plain; charset=us-ascii
8 Content-Transfer-Encoding: 7bit
9 X-From: John J Lavorato
10 X-To: Rick Buy, Ted Murphy
11 X-cc:
12 X-bcc:
13 X-Folder: \John_Lavorato_Oct2001\Notes Folders\All documents
14 X-Origin: LAVORATO-J
15 X-FileName: jlavora.nsf
16
17 How's that for VAR management.
18
19 Lavo
```

```
dataset > maildir > mcconnell-m > sent > 80
1 Message-ID: <10905805.1075843958489.JavaMail.evans@thyme>
2 Date: Mon, 7 May 2001 00:11:00 -0700 (PDT)
3 From: mike.mcconnell@enron.com
4 To: morten.pettersen@enron.com
5 Subject: Re: Tokyo opportunities
6 Mime-Version: 1.0
7 Content-Type: text/plain; charset=us-ascii
8 Content-Transfer-Encoding: 7bit
9 X-From: Mike McConnell
10 X-To: Morten E Pettersen
11 X-cc:
12 X-bcc:
13 X-Folder: \Mark_McConnell_June2001\Notes Folders\Sent
14 X-Origin: MCCONNELL-M
15 X-FileName: mmconn.nsf
16
17 I look forward to it.
18 m
```

分别是含有一个 management 和一个 opportunities 的邮件。

含有 management 的邮件排行高于含有 opportunities 的邮件的原因是 management 的文档频率低于 opportunity，搜索引擎认为频率越低的词含有更丰富的信息。

另外，在排行第八的结果中，opportunity 以复数形式存在，这里体现了去词根化的结果，可以将一个词的不同形式都识别出来。

- **opportunity, management, address**

执行结果：

```
Windows Powershell
PS C:\Users\hechu\git\enron-search-engine\src> python .\semantic_search.py
Please input your semantic query (words seperated with space): opportunity management address
('..\dataset\maildir\lavorato-j\deleted_items\190',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\lavorato-j\sent\574',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\lavorato-j\sent_items\48',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\lavorato-j\all_documents\698',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\carson-m\all_documents\125',) cosine similarity: 0.9885496183206107
('..\dataset\maildir\carson-m\_sent_mail\8',) cosine similarity: 0.9885496183206107
('..\dataset\maildir\carson-m\_sent_mail\118',) cosine similarity: 0.9885496183206107
('..\dataset\maildir\dasovich-j\all_documents\928',) cosine similarity: 0.9885496183206107
('..\dataset\maildir\sanders-r\all_documents\3460',) cosine similarity: 0.9885496183206107
('..\dataset\maildir\ruscitti-k\all_documents\262',) cosine similarity: 0.9885496183206107
semantic search time: 7.315481901168823 s
PS C:\Users\hechu\git\enron-search-engine\src> |
```

展现第一和第十的结果：

```
190
dataset > maildir > lavorato-j > deleted_items > 190
1 Message-ID: <12223725.1075857673358.JavaMail.evans@thyme>
2 Date: Tue, 29 May 2001 07:41:30 -0700 (PDT)
3 From: badeer@enron.com
4 To: john.lavorato@enron.com
5 Subject: LT management
6 Mime-Version: 1.0
7 Content-Type: text/plain; charset=us-ascii
8 Content-Transfer-Encoding: 7bit
9 X-From: Badeer, Robert
10 X-To: Lavorato, John </o=ENRON/ou=NA/cn=Recipients/cn=Jlavora>
11 X-cc:
12 X-bcc:
13 X-Folder: \jlavora\Deleted Items
14 X-Origin: Lavorato-J
15 X-FileName: jlavora.pst
16
17 It is Presto.
```

```
262 x
dataset > maildir > ruscitti-k > all_documents > 262
1 Message-ID: <10025303.1075857815129.JavaMail.evans@thyme>
2 Date: Tue, 21 Dec 1999 10:51:00 -0800 (PST)
3 From: richard.j.moller@marshmc.com
4 To: kevin.ruscitti@enron.com
5 Subject: Coun's Address
6 Mime-Version: 1.0
7 Content-Type: text/plain; charset=us-ascii
8 Content-Transfer-Encoding: 7bit
9 X-From: Richard Moller <Richard.J.Moller@marshmc.com>
10 X-To: Kevin Ruscitti
11 X-cc:
12 X-bcc:
13 X-Folder: \Kevin_Ruscitti_Dec2000\Notes Folders\All documents
14 X-Origin: Ruscitti-K
15 X-FileName: kruscit.nsf
16
17 9 Monroe Place
18 Cranbury, NJ 08512
```

- opportunity, management, address, offer

执行结果:

```
Windows Powershell
PS C:\Users\hechu\git\enron-search-engine\src> python .\semantic_search.py
Please input your semantic query (words seperated with space): opportunity management address offer
('..\dataset\maildir\lavorato-j\all_documents\698',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\lavorato-j\deleted_items\190',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\lavorato-j\sent_items\48',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\lavorato-j\sent\574',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\carson-m\sent\123',) cosine similarity: 0.9885496183206107
('..\dataset\maildir\carson-m\discussion_threads\83',) cosine similarity: 0.9885496183206107
('..\dataset\maildir\sanders-r\all_documents\3460',) cosine similarity: 0.9885496183206107
('..\dataset\maildir\carson-m\all_documents\61',) cosine similarity: 0.9885496183206107
('..\dataset\maildir\carson-m\_sent_mail\54',) cosine similarity: 0.9885496183206107
('..\dataset\maildir\dasovich-j\sent\77',) cosine similarity: 0.9885496183206107
semantic search time: 7.339409112930298 s
PS C:\Users\hechu\git\enron-search-engine\src> |
```

- opportunity, management, address, offer, price

执行结果:


```
Windows Powershell
PS C:\Users\hechu\git\enron-search-engine\src> python .\semantic_search.py
Please input your semantic query (words seperated with space): opportunity management address off
er price
('..\dataset\maildir\lavorato-j\sent_items\48',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\lavorato-j\all_documents\698',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\lavorato-j\deleted_items\190',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\lavorato-j\sent\574',) cosine similarity: 1.0770101925254814
('..\dataset\maildir\tholt-j\sent_mail\44',) cosine similarity: 1.0666666666666667
('..\dataset\maildir\dean-c\deleted_items\104',) cosine similarity: 1.0666666666666667
('..\dataset\maildir\love-p\all_documents\613',) cosine similarity: 1.0666666666666667
('..\dataset\maildir\campbell-l\sent_mail\149',) cosine similarity: 1.0666666666666667
('..\dataset\maildir\tholt-j\all_documents\183',) cosine similarity: 1.0666666666666667
('..\dataset\maildir\tholt-j\sent\44',) cosine similarity: 1.0666666666666667
semantic search time: 7.272511720657349 s
PS C:\Users\hechu\git\enron-search-engine\src> |
```

总结

通过本实验，我们实现了一个 bool 检索和语义检索系统，加深了对倒排索引，tf-idf 语义检索的理解。在优化系统过程中，进一步加深了对搜索系统的理解。

感谢助教和老师，提供了优质的学习资源！