
古诗词检索系统

计算机科学与技术专业 1711436 皮春莹

1. 问题描述

实现中华古诗词检索系统，对于给定的中华古诗词数据集（json 格式存储）创建复合索引（双词+位置），支持短语查询，支持与、或、非操作。在此基础上，进行 UI 展示以及对给定数据集的数据分析与古诗词翻译。

2. 实验环境

(1) Python 3.6

(2) JetBrains PyCharm Community Edition 2019.1.3

3. 实现思路

3.1 数据分析

相比于上次的邮件数据，这次诗词数据的数量和规范程度真让人感动。这次的数据集是 255000 首宋词和 58000 首唐诗，分批存储在 json 文件中，每个 json 文件中是一个列表，每个列表中有 1000 个元素，每个元素是一个字典，包含着一首诗词的 strains、author、paragraphs 和 title 信息。

文件名也很友好，都编好了号，所以在对诗词进行编号时，对于某个列表中的第 i 项 ($0 \leq i < 999$)，它的 ID 对应为文件名中的数字加上 i 。编号与诗词一一对应。在查询过程中，诗词的存储路径也可以通过诗词编号计算出来，不用再存储这个对应关系。

3.2 数据处理

对数据的预处理（在 `rf()` 函数中）主要是使用 `jieba` 分词，统计每首诗词中的单词。具体过程：利用文件名和数组下标对诗词进行编号（记为 `id`），每首诗词的单词信息用列表存储（记为 `single_poem`），列表的每一项是单词和位置偏移构成的元组（`word, pos`）。对于当前文件中的每一首诗词，调用 `jieba` 分词包，使用搜索引擎模式对每一句进行分词，最开始因为考虑到用户输入的单词可能不是完整的短语，在精确模式的基础上，对长词再次切分，可以提高召回率，所以选用了“搜索引擎模式”而不是“精确模式”或“全模式”。后面在实际搜索时，发现“搜索引擎模式”对于单词偏移量的确定带来了一些问题，所以最后还是改

成了使用精确模式。

分词的结果过滤掉符号等无用信息，然后将单词和该词在诗词中的偏移量存入 `single_poem` 列表，最终将文件中的所有诗词信息放入一个列表中返回。另外，为了便于完成附加功能里的数据统计分析，在这里同时统计下出现过的单词和诗人。

3.3 构建索引

在这次作业中需要建两类索引——双词索引和带位置信息的索引。双词索引是将每两个连续的词组成词对(作为短语)来索引，索引构建时将每个词对看成一个词项放到词典中。这样的话，两个词组成的短语查询就能直接处理。位置索引是在倒排记录表中，对每个 `term` 在每篇文档中的每个位置(偏移或者单词序号)进行存储。

● 索引的结构

双词索引（记为 `BiWordIndex`）是一个字典结构，字典中的每一项中，`key` 为词项，`value` 为词项所在的诗词 ID 的列表。如图所示。位置索引（记为 `PosIndex`）也是一个字典结构，字典中的每一项，`key` 为词项，`value` 是一个列表，列表中的每个元素表示一首诗词关于该词项的位置信息，仍然是一个列表来存储，0 号元素存放诗词 ID，后续的元素是 `key` 在该诗词中的偏移量。

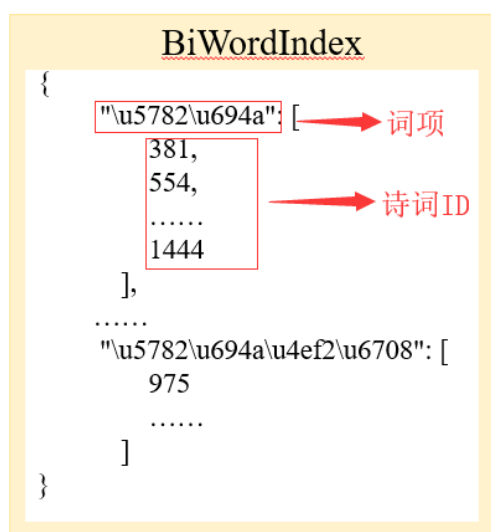


图 1 BiWordIndex

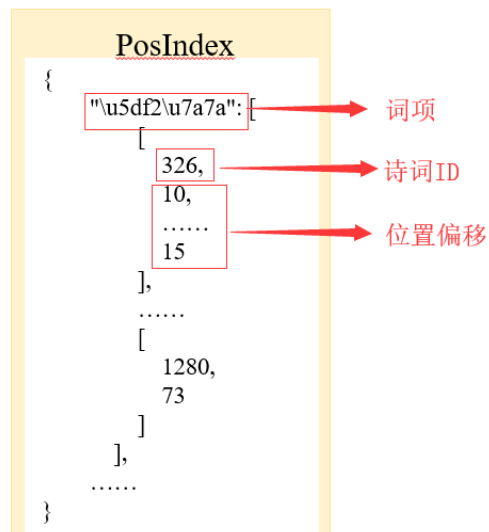


图 2 PosIndex

另外，统计诗人和单词的信息比较简单（分别记为 `Poets` 和 `Bag`），用两个列表分别存储就可以。

● 向索引中增加信息

对于每一个文件预处理得到的结果，依次处理单词集合的列表，添加双词索引时，将每两个连续的词作为一个词项，会有三种情况：

- (1) 该词项在双词索引中尚未出现，那么，加入 `{词项, [PoemID]}` 的键值对；
- (2) 该词项已在双词索引中出现，但词项对应的列表中没有当前的 `PoemID`，那么，在该列表中加入 `PoemID`；
- (3) 该词项已在双词索引中出现，并且词项对应的列表中含有当前的 `PoemID`，

那么，直接处理下一个词项。

利用元组（单词 `term`，偏移 `pos`）向位置索引 `PosIndex` 添加信息时，会有以下三种情况：

- (1) `term` 不属于 `PosIndex` 的关键字，那么，加入 `key` 为 `term`，`value` 为 `[[PoemID,pos]]` 的键值对；
- (2) `term` 是已经存在的关键字，但它对应的列表中未出现 `PoemID`，那么，在他现有列表中附加一个列表 `[PoemID, pos]`；
- (3) `term` 是已经存在的关键字，并且对应的列表中的某个内层列表的 0 号元素是 `PoemID`，那么，在内层列表中追加元素 `pos`。

最后，将索引写入磁盘。由于位置索引的文件体积太大，所以在实际编程的过程中，我将唐诗宋词分成了三个文件夹，分块建了索引，每次只加载一部分信息，有需要时，才加载更多的索引。

3.4 数据查询

从用户图形化界面中获取到的是一个字符串，指明需要进行的短语查询、与或非操作，查找索引，返回相关诗词的编号，利用编号映射到诗词所在的文件，读出诗词内容。

● 提取查询信息

我设计的输入的字符串格式为 “`xxxx or/and/andnot xxxx or/and/andnot xxxx……`”，其中，`xxxx` 是单词、短语或诗句，由 `and`、`or` 或 `andnot` 相连接，表示这些短语之间的与或非关系。在这里，考虑到优先级的关系，提取查询信息主要分为以下三个步骤：

- (1) 先利用 “`or`” 拆分字符串，把拆分结果记录在 `or_terms` 列表中。
- (2) 对于其中的每一项，继续用 “`and`” 拆分，若该词项中不含 “`and`”，则不做操作；若该词项中含有 “`and`”，则将拆分的结果记录在 `and_terms` 列表中，并且从 `or_terms` 中删去该词项。
- (3) 对于 `and_terms` 中的每一项，用 “`not`” 进行拆分，若该词项中不含 “`not`”，则不做操作；若该词项中含有 “`not`”，则将拆分的结果记录在 `not_terms` 列表中，并且从 `and_terms` 中删去该词项。

至此，`or_terms`、`and_terms`、`not_terms` 中各自分别存着需要与其他两个列表进行或操作、与操作、非操作的词项；而在 `or_terms` 列表的内部，词项之间是或的关系，在 `and_terms` 和 `not_terms` 的内部，词项之间是与的关系。而对于每一个列表中的每一个词项，若词项是一个单词，那么不用做过多的处理，直接查索引，提取诗词编号的列表；如果词项不仅仅是一个单词，而是两三个单词组成的短语，或者是一句诗，那么在查找时需要保证单词之间的位置关系，首先要对该词项进行分词，然后将每两个相邻的单词组合成一个词项，在索引中找到该次相对应的诗词编号列表，多个列表之间进行 “与” 操作。

由于双词索引体积小，查询速度较快，但对查询词项的要求较高，有的词项在双词索引中可能搜索结果为空；而位置索引记录的信息较为全面，能提高召回率，但体积过大，搜索速度慢。因此在查索引时，利用双词索引对计算代价比较高的短语进行索引，若返回结果不够满意，再利用位置索引对短语进行处理，这样就充分发挥各自的优势。

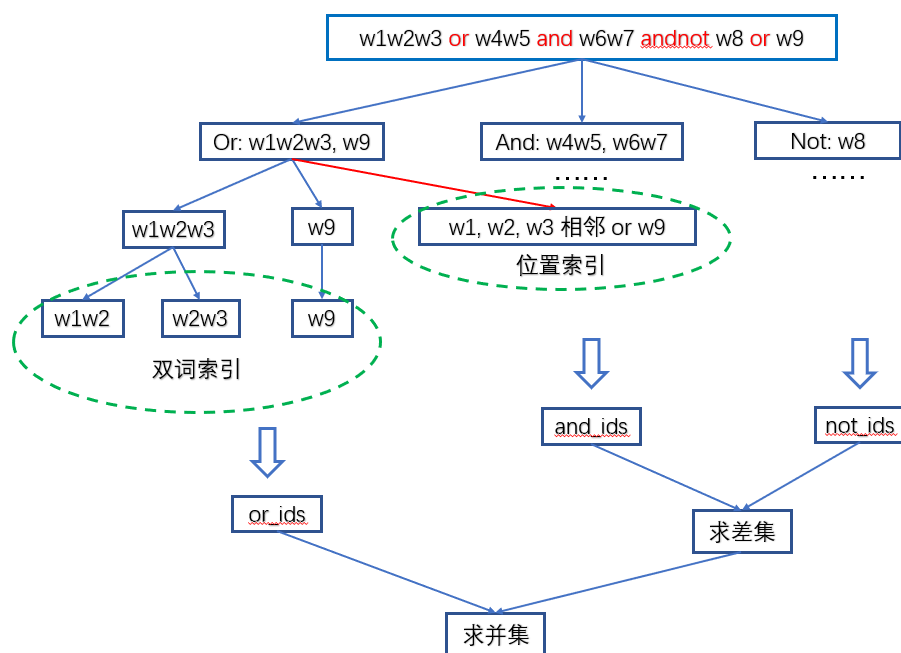


图 3 诗词查询

在用位置索引查询时，对于同一种运算关系中的由短语分词得到的的单词列表，每次取一个单词，用字典存储它涉及到的诗词 ID 及对应的偏移信息，以文章 id 为主键，值为列表，列表中存储关键词出现过的位置。以第一个单词为基准，每次将当前第 m 个单词的字典与它的字典进行对比，若在两个字典中，对于同一篇文章，得到的两个列表中出现位置相差为 m，证明当前 id 符合条件。

以“明月清風 or 平世功名”查询为例，先用“or”“and”“not”依次拆分，得到的结果：or_terms=['明月清風', '平世功名'], and_terms=[], not_terms=[]。对 or_terms 中的每一项进行分词，分别得到 ['明月', '清風']和['平世', '功名']。

先从磁盘加载双词索引，将分词结果的每两项一组合，两个列表分别得到了词项：“明月清風”、“平世功名”，这里不同于最初 or_terms 中元素的是，or_terms 中的元素可能是单词或两个及两个以上的单词组成的短语，但这里的处理结果一定是单词或两个单词组成的短语。

组合成词项之后，在双词索引中查找每个词项，记录他们对应的列表，多个列表分别排序去重后，进行或操作，得到对于 or_terms 的诗词编号列表 or_ids[]，由于 and_terms 和 not_terms 为空，所以对应的 and_ids[]和 not_ids[]也为空，将 or_ids[]和 and_ids[]进行或操作，再与 not_ids[]进行非操作，得到由双词索引检索的结果。判断首次检索的结果是否满意，这里我的判断标准是搜索到的诗词数量，如果数量较少的话，需要继续去位置索引中查找。

如果上面的例子需要在位置索引中进行查找，仍然需要使用逻辑词拆分、结巴分词之后的结果：['明月', '清風']和['平世', '功名']。先查找“明月”对应的列表，若关键词不存在，则返回空列表，再查找“清風”对应的列表，若关键词不存在，则返回空列表。当上面两个列表都不为空时，例如找到的分别是 [[12,5,10],[20,8,20],[25,10]], [[28,7],[12,9],[9,17,30]]，找出每个涉及到的诗词，即内层列表的首个元素，使用字典结构，以诗词 id 为主键，值为列表，列表中存储关键词出现过的位置，若在列表中出现连续长度等于外层列表数量时，证明当前 id 符合条件。在这里形成的字典是 {9: [17,30], 12: [5,9,10], 20: [8,20], 25: [10], 18:

[7]}, 检查可知, id=12 的诗词满足条件, 对['平世','功名']也进行相同操作, 将得到的结果进行或操作, 返回诗词列表。

3.5 附加功能

在完成主要的功能之外, 还编写了用户图形化界面。这个主要使用了 python 的窗口视窗设计的模块 tkinter, 放置了用户输入框和结果展示的文本框。在展示查询结果时, 由上面的介绍可知, 搜索功能的函数返回的是诗词 ID 集合, 需要在这里将 ID 翻译成对应的文件名, 具体的对应关系在“数据分析”中已经说明。这里设计的每次只展示一首诗词的名称、作者和内容, 文本框下面放置了“Next”按钮, 点击按钮时, 系统会根据下一个诗词 ID, 加载诗词内容并显示在文本框中。

界面的左侧是关键字排名分析和作者产量分析。由于在数据处理的过程中, 已经将分词产生的所有单词存入了 Bag, 将诗人的名字存入了 Poets, 这里只需要对数据进行汇总和展示。我计划将关键字和作者的信息做成多张词云图, 然后将这多张图片合成 GIF 动图, 在界面右侧展示出来。词云是对文本数据中出现频率较高的“关键词”在视觉上的突出呈现, 在 python 中可以用 wordcloud 包比较轻松地实现: 加载 Bag 和 Poets, 分别将两个列表内的元素用空格连接, 然后设置词云的背景色、字号、字体、配色等属性, 这里需要注意的是, wordcloud 包不支持繁体中文, 所以我另外下载了方正黑体繁体的字体包并进行了设置 (font_path="FangZhengHeiTiFanTi-1.ttf"), 之后将生成的词云图保存成 PNG 格式, 并将多张 PNG 图片合成 GIF。

由于没有找到 python 的 tkinter 包中可以直接展示 GIF 的控件, 我编写了 ImageLabel 类, 用于展示 GIF 图。ImageLabel 类中定义了 load()、unload()和 next_frame()函数, load()使用了 PIL 的 Image 和 ImageTk 工具包, 先调用 Image 的 load()函数加载 GIF, 每隔固定的时间段, 循环展示一帧, 循环展示的过程由 next_frame()控制。界面如下图所示:



图 4 初始界面

4. 代码实现

本实验分为三个文件：splitSentence.py、draw_graph.py 和 search_words.py。其中，splitSentence.py 是对数据进行处理，并建立双词索引和位置索引。draw_graph.py 是将关键词和作者信息做成词云图。search_words.py 是显示查询界面，根据用户的输入，加载索引并进行查询，返回搜索结果。

4.1 splitSentence.py

读取诗词文件，对于每首诗词，使用 jieba 的精确模式进行分词，并给每个词附加位置偏移信息，将结果存入列表。

```
1. def rf(path, obj):
2.     # 返回的应该是一个长度为 1000 的列表，每一项是当前文档内的一首古诗的内容分词结果
3.     tmp_file = open(path, 'r', encoding='utf-8')
4.     tmp_list = json.loads(tmp_file.read())
5.     for j in range(1000):
6.         # 对于每首诗词，先提取正文部分
7.         sentence_list = tmp_list[j]["paragraphs"]
8.         poet = tmp_list[j]["author"]
9.         poem_index.poets.append(poet) # 便于完成附加功能里的数据统计分析
10.        word_count = 0
11.        # 每首诗词的单词信息用列表存储，列表的每一项是单词和位置偏移构成的元组
12.        single_poem = []
13.        for sentence in sentence_list:
14.            # 使用 jieba 的精确模式进行分词，把结果放入 obj
15.            words = jieba.lcut(sentence)
16.            for word in words:
17.                # 过滤掉符号等无用信息
18.                if word not in [' ', '。', '!', '?', ':', '{', '}', '/', '(', ')']:
19.                    # 给每个词附加位置偏移信息
20.                    single_poem.append((word, word_count))
21.                    word_count += 1
22.            obj.append(single_poem)
23.        tmp_file.close()
```

定义了 Index 类，便于整理所有的数据处理结果、双词索引、位置索引，Index 类中定义了 add_index(self, poem_id, word_tuples) 函数，参数 poem_id 是根据文件名和诗词的下标计算得到的，计算 poem_id 的代码比较简单，不再粘贴。word_tuples 是经过上面的 rf() 函数处理过的文件中的一首诗词（每个文件中有一千首），add_index() 的流程图如下：

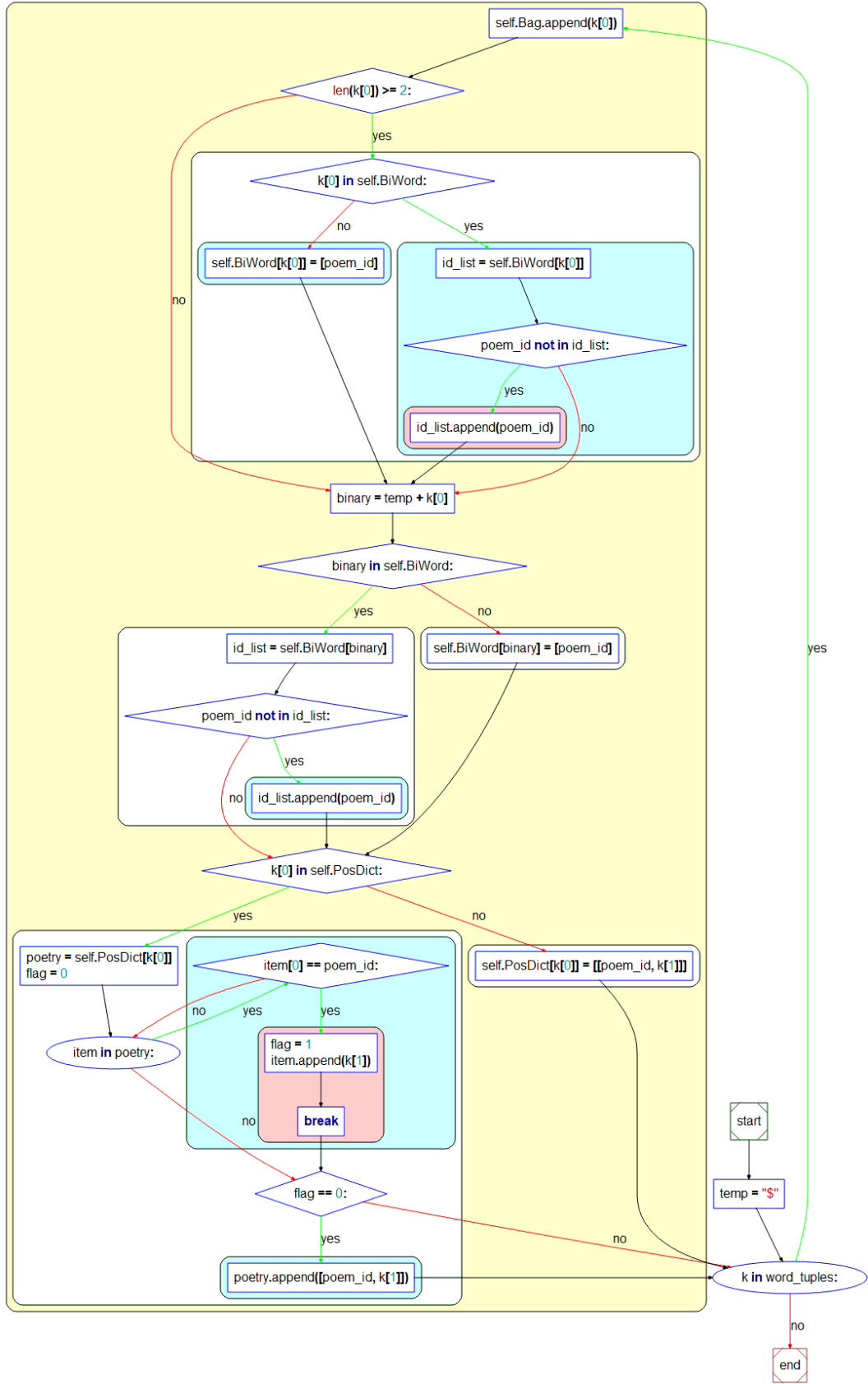


图 5 ClusterControlFlow-add_index

下面是 Index 类的定义：

```
1. class Index():
2.     def __init__(self):
3.         self.PosDict = dict() # 位置索引
4.         self.BiWord = dict() # 双词索引
5.         # 统计诗人和单词的信息
6.         self.Bag = list()
7.         self.poets = list()
8.
9.         # 将每篇文章得到的单词集合插入索引
10.    def add_index(self, poem_id, word_tuples):
11.        temp = "$"
12.        # words 是一个元素为元组的列表
13.        for k in word_tuples:
14.            # k=(word,pos)
15.            self.Bag.append(k[0])
16.            if len(k[0]) >= 2:
17.                if k[0] in self.BiWord:
18.                    id_list = self.BiWord[k[0]]
19.                    if poem_id not in id_list:
20.                        id_list.append(poem_id)
21.                else:
22.                    self.BiWord[k[0]] = [poem_id]
23.            # 添加双词索引，
24.            # 将每两个连续的词组成一个词项放到词典中
25.            binary = temp + k[0]
26.            # 会有三种情况：
27.            # (1)词项在双词索引中尚未出现，加入{词项，[PoemID]}的键值对；
28.            # (2)词项已在双词索引中出现，但词项对应的列表中没有当前的 PoemID，在该
            # 列表中加入 PoemID；
29.            # (3)词项已在双词索引中出现，并且词项对应的列表中含有当前的 PoemID，直
            # 接处理下一个词项。
30.            if binary in self.BiWord:
31.                id_list = self.BiWord[binary]
32.                if poem_id not in id_list:
33.                    id_list.append(poem_id)
34.            else:
35.                self.BiWord[binary] = [poem_id]
36.
37.            # 添加位置索引
38.            # key 为词项，value 是一个双层的列表，内层列表表示一首诗词关于该词项的
            # 位置信息，
39.            # 0 号元素存放诗词 ID，后续的元素是 key 在该诗词中的偏移量
40.            if k[0] in self.PosDict:
```



```

41.         # term 是已经存在的关键字
42.         poetry = self.PosDict[k[0]]
43.         # 在 poetry 列表中找编号 id 的列表，并加入位置信息
44.         flag = 0 # 0 表示 poem_id 不在当前的列表中
45.         for item in poetry:
46.             if item[0] == poem_id:
47.                 # 对应的列表中的某个内层列表的 0 号元素是 PoemID，那么，在内层列表
                    中追加元素 pos
48.                 flag = 1
49.                 item.append(k[1])
50.                 break
51.         if flag == 0:
52.             # 对应的列表中未出现 PoemID，那么，在现有列表中附加一个列表
                [PoemID, pos]
53.             poetry.append([poem_id, k[1]])
54.         else:
55.             # term 不属于 PosIndex 的关键字，那么，加入 key 为 term，value 为
                [[PoemID,pos]]的键值对；
56.             self.PosDict[k[0]] = [[poem_id, k[1]]]

```

剩余的代码是递归读取文件夹下的每个文件，以及将索引写入磁盘，比较简单，不再粘贴在这里。

4.2 draw_graph.py

加载 Bag 和 Poets，用 wordcloud 包制作词云图：

```

1. from wordcloud import WordCloud
2. import matplotlib.pyplot as plt
3. import json
4.
5.
6. def draw_graph():
7.     with open("song_words_bag.json") as f:
8.         mytext1 = json.load(f)
9.     f.close()
10.    w1 = " ".join(mytext1)
11.    # 加载背景图片
12.    # cloud_mask = np.array(Image.open("cloud.png"))
13.    # 设置词云
14.    wc = WordCloud(background_color="white", # 设置背景颜色
15.                   # mask=cloud_mask, # 设置背景图片
16.                   max_words=2000, # 设置最大显示的字数
17.                   # stopwords = "", #设置停用词

```

```

18.         font_path="FangZhengHeiTiFanTi-1.ttf",
19.         # 设置中文字体，使得词云可以显示（词云默认字体是
           "DroidSansMono.ttf 字体库”，不支持中文）
20.         max_font_size=50, # 设置字体最大值
21.         random_state=30, # 设置有多少种随机生成状态，即有多少种配色
           方案
22.     )
23.     myword = wc.generate(wl) # 生成词云
24.
25.     # 展示词云图
26.     plt.imshow(myword)
27.     plt.axis("off")
28.     plt.show()
29.     # 存储图片
30.     wc.to_file('a4.png')
31.
32. if __name__ == '__main__':
33.     draw_graph()

```

4.3 search_words.py

界面上的控件比较简单，这里不再解释相关的代码。按键“Search”绑定了 hit_btn()函数，用于读取用户输入的字符串，之后调用 search_poem()函数，得到相关诗词的编号，并将搜索结果中的第一首诗词信息展示在文本框中。

```

1. def hit_btn():
2.     # 定义一个函数功能供点击 Button 按键时调用，调用命令参数 command=函数名
3.     # 搜索结果总数
4.     global total
5.     # 搜索结果
6.     global ans_ids
7.     # 当前展示的诗词的计数
8.     global tmp_num
9.
10.    # 获取用户输入的信息
11.    e_from = var_from.get()
12.    print(e_from)
13.    ans_ids = search_poem(e_from, "song")
14.    total = len(ans_ids)
15.    tmp_num = 0
16.    if total > 0:
17.        t.delete(0.0, END)
18.        # 对应到文件名
19.        file_id = int(ans_ids[0] / 1000) * 1000

```

```

20.         print("The first fileID and poemID-in-file to be show is")
21.         print(file_id)
22.         num = ans_ids[0] - file_id
23.         print(num)
24.         # 由于唐诗和宋词是分块建的索引，所以要先判断需要加载哪个文件
25.         if ans_ids[0] < 255000:
26.             # 宋词
27.             path = "E:\\ir\\chinese poetry\\poet.song." + str(file_id) + ".json"
28.         else:
29.             # 唐诗
30.             path = "E:\\ir\\chinese poetry\\poet.tang." + str(file_id) + ".json"
31.
32.         with open(path, 'r', encoding='UTF-8') as f:
33.             tp = json.load(f)
34.             f.close()
35.             # 将搜索结果中的第一首诗词信息展示在文本框中
36.             poem_dict = tp[num]
37.             author = poem_dict["author"]
38.             para = poem_dict["paragraphs"]
39.             title = poem_dict["title"]
40.             t.insert(INSERT, title + "\n")
41.             t.insert(INSERT, author + "\n")
42.             for i in para:
43.                 t.insert(INSERT, i + "\n")
44.         else:
45.             t.delete(0.0, END)
46.             t.insert(INSERT, "未搜索到相关内容\n")

```

“Next”按钮绑定了 `show_btn()` 函数，通过改变全局变量 `tmp_num` 的值，来切换文本框的内容，代码比较简单，不再粘贴。

在 `hit_btn()` 中调用了 `search_poem()` 函数，`search_poem()` 的作用是根据用户输入，返回相应的诗词 ID 集合。在 `search_poem()` 中，根据逻辑运算的优先级，将查询语句拆分成三个列表（`or_terms`、`and_terms` 和 `not_terms`），具体的拆分方法在实验思路中已经阐述。然后加载双词索引。接下来分别对每个列表进行操作，列表内部的每个词项，是一个或多个相邻的短语，处理的方法类似，写在 `internal_term()` 中，在这里直接对每个词项调用该函数。对于列表中每个词项得到的 PoemID 列表，在 `or_terms` 中需要进行或操作，因为只要有一个满足就可以选择该 ID，在 `not_terms` 中也需要进行或操作，因为只要有一个不满足就可以过滤该 ID，但在 `and_terms` 中需要进行与操作，因为需要满足全部词项的 ID 才可以被留下，最后调用 `merge_vector()` 将 `or_terms`、`and_terms` 和 `not_terms` 得到的诗词编号列表进行合并。下面是 `search_poem()` 的前半部分，是关于双词索引的部分：

```

1. or_terms = []

```

```
2. and_terms = []
3. not_terms = []
4. # 考虑优先级的关系，提取查询信息
5. # 先利用“or”拆分字符串，把拆分结果记录在 or_terms 列表中
6. or_terms.extend(quiry.split('or'))
7. for each in or_terms:
8.     print(each)
9.     and_terms.extend(each.split('and'))
10. # 对于其中的每一项，继续用“and”拆分，若该词项中不含“and”，则不做操作；
11. if 'and' in each:
12.     # 若该词项中含有“and”，则将拆分的结果记录在 and_terms 列表中，并且从
    or_terms 中删去该词项
13.     or_terms.remove(each)
14. for each in and_terms:
15.     # 对于 and_terms 中的每一项，用“not”进行拆分，若该词项中不含“not”，则不做操
    作
16.     if 'not' in each:
17.         # 若该词项中含有“not”，则将拆分的结果记录在 not_terms 列表中，并且从
    and_terms 中删去该词项
18.         not_terms.append(each.replace('not', ''))
19.         and_terms.remove(each)
20.
21. print(or_terms)
22. print(and_terms)
23. print(not_terms)
24.
25. filepath = prefix + "_poem_BiWordIndex.json"
26. # 先进行双词索引
27. with open(filepath, 'r', encoding='UTF-8') as f:
28.     BiWdIndex = json.load(f)
29. f.close()
30. # 要保证插入的 id 是有序且不重复的，集合（set）是一个无序的不重复元素序列
31. not_ids = []
32. ni = 0
33. nm = [[], [], []]
34. for each in not_terms:
35.     nm[ni] = internal_term(each, BiWdIndex)
36.     not_ids.extend(nm[ni])
37.     ni += 1
38.
39. tm = [[], [], []]
40. ai = 0
41. # 两层求交集操作
42. for each in and_terms:
```

```

43.     tm[ai] = internal_term(each, BiWdIndex)
44.     ai += 1
45. and_ids = tm[0]
46. ii = 1
47. while ii < 3 and len(tm[ii]) > 0:
48.     for i in reversed(and_ids):
49.         # 反向删除
50.         if i not in tm[ii]:
51.             and_ids.remove(i)
52.     ii += 1
53.
54. or_ids = []
55. oi = 0
56. om = [[], [], []]
57. for each in or_terms:
58.     om[oi] = internal_term(each, BiWdIndex)
59.     or_ids.extend(om[oi])
60.     oi += 1
61.
62. poemIDs = merge_vector(not_ids, and_ids, or_ids)
63. count = len(poemIDs)
64. print("poemIDs: ")
65. print(poemIDs)

```

如果双词索引中未能找到符合条件的诗词，要继续在位置索引中查找。下面是 `search_poem()` 的后半部分，主要是关于位置索引的查找。每个列表中的每个短语，需要先进行分词，对于分词得到的若干个单词，处理方法类似，写在 `merge_with_pos()` 中，这里直接调用，得到的是该短语对应的是此编号列表，然后与上面类似，`or_terms`、`and_terms`、`not_terms` 的词项间采取相应的操作，最后调用 `merge_vector()` 将 `or_terms`、`and_terms` 和 `not_terms` 得到的诗词编号列表进行合并。

```

1. # 如果需要的话再进行位置索引
2. if count == 0:
3.     # if flag == 1:
4.     print("----In PosIndex----")
5.     filepath2 = prefix + "_poem_PosIndex.json"
6.     with open(filepath2, 'r', encoding='UTF-8') as f:
7.         PosIndex = json.load(f)
8.     f.close()
9.
10.    not_ids2 = []
11.    for each in not_terms:
12.        ws = jieba.lcut(each)
13.        not_ids2.extend(merge_with_pos(ws, PosIndex))

```

```

14.
15.     and_ids2 = []
16.     cn = len(and_terms)
17.     if cn > 0:
18.         ws = jieba.lcut(and_terms[0])
19.         and_ids2.extend(merge_with_pos(ws, PosIndex))
20.         ind = 1
21.         while ind < cn:
22.             and_ids3 = []
23.             ws = jieba.lcut(and_terms[ind])
24.             and_ids3.extend(merge_with_pos(ws, PosIndex))
25.             for i in reversed(and_ids2):
26.                 # 反向删除
27.                 if i not in and_ids3:
28.                     and_ids2.remove(i)
29.             ind += 1
30.
31.     or_ids2 = []
32.     for each in or_terms:
33.         ws = jieba.lcut(each)
34.         or_ids2.extend(merge_with_pos(ws, PosIndex))
35.
36.     return merge_vector(not_ids2, and_ids2, or_ids2)
37. else:
38.     return poemIDs

```

接下来对 `search_poem()` 中调用的 `internal_term()` 进行说明。`internal_term()` 是对参数中的短语进行分词，然后将每两个相邻的单词组合成一个词项，对每个词项的查询结果取交集，返回编号列表。

```

1. def internal_term(each, BiWdIndex):
2.     ws = jieba.lcut(each)
3.     tmp_ids = BiWdIndex.get(ws[0], [])
4.     tm = tmp_ids
5.     if len(ws) > 1:
6.         # 对该词项进行分词，然后将每两个相邻的单词组合成一个词项
7.         for ww in range(len(ws) - 1):
8.             tmp_ids = []
9.             obj = ws[ww] + ws[ww + 1] # 构造双词结构
10.            tmp_ids.extend(BiWdIndex.get(obj, []))
11.
12.            ss = set(tmp_ids)
13.            tmp_ids = list(ss)
14.            tmp_ids.sort()
15.            # 对每个词的结果要取交集

```

```
16.         # 如果不在当前 tmp_ids 中出现, 则删去
17.         if len(ws) == 2:
18.             return tmp_ids
19.         for tid in reversed(tm):
20.             # 反向删除
21.             if tid not in tmp_ids:
22.                 tm.remove(tid)
23.     return tm
```

merge_vector()的参数是三个数组, 分别是由前面逻辑词拆分得到的词项的查询结果, 记为 not_ids、and_ids 和 or_ids。这里需要对列表进行去重和排序。之后, 在 and_ids 中去掉 not_ids 所含的 ID, 再与 or_ids 求并集, 关键代码如下:

```
1. len2 = len(not_list)
2. len1 = len(and_list)
3. res = []
4. p1 = 0
5. p2 = 0
6. if len2 == 0 and len1 == 0:
7.     print("仅 or 短语组合")
8.     return or_list
9. # 求差集
10. while p1 < len1 and p2 < len2:
11.     while p1 < len1 and p2 < len2 and and_list[p1] < not_list[p2]:
12.         res.append(and_list[p1])
13.         p1 += 1
14.     while p2 < len2 and p1 < len1 and and_list[p1] > not_list[p2]:
15.         p2 += 1
16.     if p2 < len2 and p1 < len1 and and_list[p1] == not_list[p2]:
17.         p1 += 1
18.         p2 += 1
19. while p1 < len1:
20.     res.append(and_list[p1])
21.     p1 += 1
22.
23. len3 = len(res)
24. len4 = len(or_list)
25. p3 = 0
26. p4 = 0
27. res_list = []
28. # 列表求并集
29. while p3 < len3 and p4 < len4:
30.     while res[p3] < or_list[p4] and p3 < len3:
31.         res_list.append(res[p3])
32.         p3 += 1
```

```

33.     while res[p3] > or_list[p4] and p4 < len4:
34.         res_list.append(or_list[p4])
35.         p4 += 1
36.     if res[p3] == or_list[p4]:
37.         res_list.append(res[p3])
38.         p3 += 1
39.         p4 += 1
40. if p3 == len3:
41.     while p4 < len4:
42.         res_list.append(or_list[p4])
43.         p4 += 1
44. else:
45.     while p3 < len3:
46.         res_list.append(res[p3])
47.         p3 += 1

```

用位置索引查询时，调用了 `merge_with_pos()` 函数，它的功能主要是筛选同一首诗词中单词出现偏移位置满足条件的诗词 ID。具体做法是：对于参数传进来的单词列表，每次取一个单词，用字典存储它涉及到的诗词 ID 及对应的偏移信息，以文章 id 为主键，值为列表，列表中存储关键词出现过的位置。以第一个单词为基准，每次将当前第 `m` 个单词的字典与它的字典进行对比，若在两个字典中，对于同一篇文章，得到的两个列表中出现位置相差为 `m`，证明当前 id 符合条件，代码的关键部分如下：

```

1. tmp_dict = {}
2. p = index.get(words[0], [])
3. print("for each word, take out the poemID-postion-list : ")
4. for li in p:
5.     print(li)
6.     if li[0] not in tmp_dict:
7.         tmp_dict[li[0]] = []
8.         for i in range(1, len(li)):
9.             tmp_dict[li[0]].append(li[i])
10. pids = list(tmp_dict.keys())
11. if c_word >= 2:
12.     m = 1
13.     while m < c_word:
14.         p = index.get(words[m], [])
15.         tmp_dict1 = {}
16.         # print("for each word, take out the poemID-postion-list : ")
17.         for li in p:
18.             # print(li)
19.             if li[0] not in tmp_dict1:
20.                 tmp_dict1[li[0]] = []
21.                 for i in range(1, len(li)):

```



```

22.         tmp_dict1[li[0]].append(li[i])
23.     t2 = list(tmp_dict1.keys())
24.     for i in reversed(pids):
25.         # 反向删除
26.         if i not in t2:
27.             pids.remove(i)
28.     for pid in reversed(pids):
29.         tmp_list = tmp_dict[pid]
30.         tmp_list.sort()
31.         tmp_list1 = tmp_dict1[pid]
32.         tmp_list1.sort()
33.         j = 0
34.         flag = 0
35.         for tid in tmp_list:
36.             while j < len(tmp_list1):
37.                 if tid == tmp_list1[j] - m: # m 表示第 m 个词，同时也是与第
一个词的间距
38.                     flag = 1
39.                     break
40.                 j += 1
41.                 if flag == 1: # PoemID 得到了确认
42.                     break
43.                 if flag == 0:
44.                     pids.remove(pid)
45.             m += 1

```

至此，查询功能已经完成，系统可以满足短语、与或非关系混合的查询。为了展示对关键词和世人的分析结果，我将 `draw_graph.py` 文件生成的 PNG 格式的词云图合成了 GIF，然后编写了 `ImageLabel` 类，用于显示 GIF，`ImageLabel` 类显示 GIF 的主要原理是获取 GIF 的所有的帧，每个一段时间，显示下一帧，不断循环。代码如下：

```

1. class ImageLabel(tk.Label):
2.     def load(self, im):
3.         if isinstance(im, str):
4.             im = Image.open(im)
5.             self.loc = 0
6.             self.frames = []
7.
8.         try:
9.             for i in count(1):
10.                 self.frames.append(ImageTk.PhotoImage(im.copy()))
11.                 im.seek(i)
12.             except EOFError:
13.                 pass

```

```

14.
15.     try:
16.         self.delay = im.info['duration']
17.     except:
18.         self.delay = 250
19.
20.     if len(self.frames) == 1:
21.         self.config(image=self.frames[0])
22.     else:
23.         self.next_frame()
24.
25.     def unload(self):
26.         self.config(image=None)
27.         self.frames = None
28.
29.     def next_frame(self):
30.         if self.frames:
31.             self.loc += 1
32.             self.loc %= len(self.frames)
33.             self.config(image=self.frames[self.loc])
34.             self.after(self.delay, self.next_frame)

```

5. 结果展示

本实验中的诗词检索系统支持短语、与或非关系混合的查询，例如：“明月清風 or 平世功名”、“流水 and 行舟 andnot 古渡頭”、“西風立馬頻回首”等。

(1) 输入“古渡頭 or 行舟”，点击“Go”：



图 6 古渡頭 or 行舟 (1)

点击“Next”:



图 7 古渡頭 or 行舟 (2)

下面是在这个过程中，控制台的输出信息：

```
古渡頭or行舟
古渡頭
行舟
['古渡頭', '行舟']
['古渡頭', '行舟']
[]
Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\asus\AppData\Local\Temp\jieba.cache
in merge_vector(): not_ids, and_ids, or_ids:
Loading model cost 0.870 seconds.
Prefix dict has been built successfully.
[]
[24122]
[429, 576, 2164, 2610, 6227, 7052, 8366, 9800, 11185, 11199, 11308, 11493, 11889,
poemIDs:
[429, 576, 2164, 2610, 6227, 7052, 8366, 9800, 11185, 11199, 11308, 11493, 11889,
The first fileID and poemID-in-file to be show is
0
429
clicked show-btn , total num and tempID :
149
576
clicked show-btn , total num and tempID :
149
```

图 8 控制台的输出信息

(2) 输入“流水 and 行舟 andnot 古渡頭”，点击“Go”:



图 11 明月清風 or 平世功名(1)

点击“Next”:

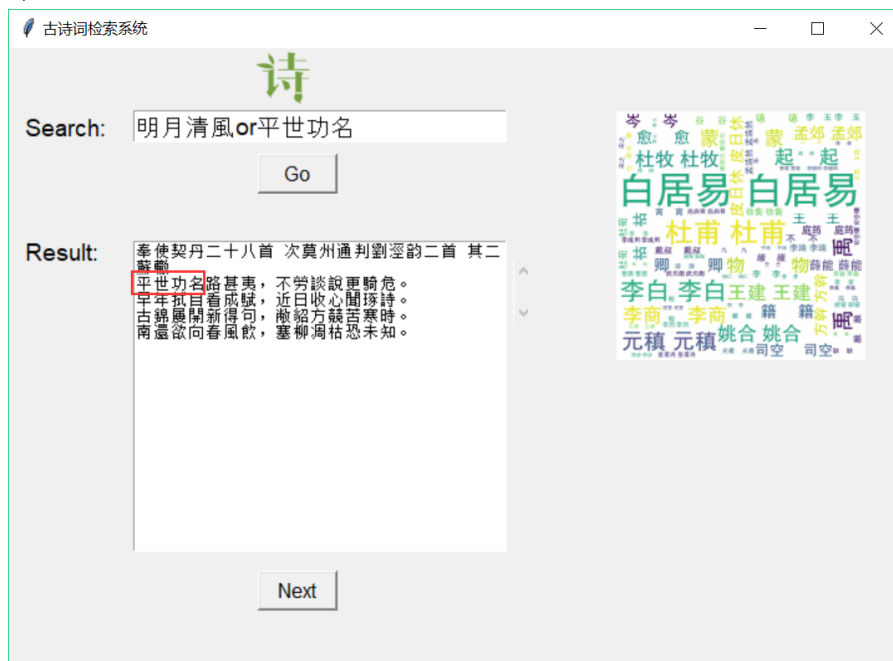


图 12 明月清風 or 平世功名(2)

(5) 输入“芭蕉雨 and 梨花地 and 梁州”，点击“Go”:



图 13 芭蕉雨 and 梨花地 and 梁州

还进行了一些其他的查询测试，记录在表 1 中，不再给出图片。

表 1 查询测试

输入	控制台输出结果
古渡頭 or 行舟	178, 429, 576, 2164, 2610, 2936, 6227……
行舟 andnot 古渡頭	429, 576, 2164, 2610, 6227, 7052, 8366……
流水 and 行舟 andnot 古渡頭	51972
流水 or 行舟 or 古渡頭	25, 37, 167, 178, 196, 246, 429, 576……
春風 or 流水 or 楊柳	25, 37, 167, 178, 179, 196, 218, 223……
古往今來	1713, 1817, 3553, 11886, 12942, 31547……
月 and 水	141, 240, 278, 962, 1713, 1770, 2087……
明月清風	1811, 5687, 17888, 19708, 20052, 26432……
平世功名	10825, 13926, 48936, 62139, 62359, 62499……
明月清風 or 平世功名	405, 869, 1811, 5687, 7248, 9661, 10825……
明月清風 and 平世功名	未搜索到相关内容

6. 实验总结

优点

- (1) 支持短语查询、逻辑查询以及短语和逻辑混合的查询
- (2) 查询结果的召回率较高
- (3) 查询结果的正确率很高

缺点

- (1) 程序的鲁棒性不够好，没有做过多的错误捕获和处理，输入的查询字符串需要符合规范，才能正常地进行查询。
- (2) 程序有点冗长，查询时间较长，由于用户输入的短语基本都是由两三个短语组成，所以一些地方还可以继续简化。可以换用更高效的数据结构。