

# 信息检索系统实验1

1911590 周安琪

## 1 数据集

## 2 索引构建与检索 (40%)

作业的第一部分是使用**blocked sort-based indexing (BSBI)** 算法来构建倒排索引并实现布尔检索。关于BSBI算法可以参考老师课件或者斯坦福教材 [Section 4.2](#)。以下摘自教材内容

### 2.1 IdMap

把term映射到int类型的termID

所以在往map中添加的时候，肯定是以term为key。

```
1 def _get_str(self, i):
2     return self.id_to_str[i]
3
4 def _get_id(self, s):
5     # 如果参数中的str在map中没有，那么加进去
6     if s not in self.str_to_id:
7         self.str_to_id[s] = len(self.id_to_str)
8         self.id_to_str.append(s)
9     return self.str_to_id.get(s)
```

### 2.2 将倒排列表编码成字节数组

### 2.3 磁盘上的倒排索引

## 2.4 索引

### 2.4.1 解析

`parse_block` 对文件做语法分析，生成 `term-ID->docID` 的对，并且把这些对存储在内存里，直到收集一个block停止。

```
1  #参数是str类型，是要处理的文件的地址。
2  def parse_block(self, block_dir_relative):
3
4      td_pairs = []
5
6      #首先把两个路径连起来，然后把里面的文档排序，按次序打开file
7      for file_dir in sorted(os.listdir(os.path.join(self.data_dir,
8          block_dir_relative))):
9
10         # 'r': 以只读方式打开文件
11         with open(os.path.join(self.data_dir, block_dir_relative,
12             file_dir), 'r') as f:
13             # strip()删除空白符 (包括'\n', '\r', '\t', ' ')
14             # split()以空格和换行为分隔符分开，返回分割后的字符串列表。
15             content = f.read().strip().split() # list of tokens
16
17             # 把doc映射到doc_id
18             doc_id =
19             self.doc_id_map[os.path.join(block_dir_relative, file_dir)]
20
21             # 现在content是由一堆token组成的东西
22             for token in content:
23                 # 把token映射到id，然后把pair加入到列表中
24                 term_id = self.term_id_map[token]
25                 td_pairs.append([term_id, doc_id])
26
27     return td_pairs
```

### 2.4.2 倒排表

#### 2.4.2.1 append()

把做好的 `term->posting_list` 加到索引文件的末尾

- 用posting\_encoding来给postings\_list编码
- 将metadata 以 `self.terms` 和 `self.postings_dict` 的形式储存，  
`self.postings_dict` 把 `termID` 映射到一个三元组
- 加入磁盘里的索引文件

```

1  def append(self, term, postings_list):
2      # 编码
3      encoded_postings_list =
4      self.postings_encoding.encode(postings_list)
5      # 开始地点是文件末尾，目前的偏移量是0
6      start_position_in_index_file = self.index_file.seek(0, 2)
7      # 返回写入的字符长度。
8      length_in_bytes_of_postings_list = self.index_file.write(
9          encoded_postings_list)
10
11     self.terms.append(term)
12     # termID映射到一个三元组，这个三元组其实就是倒排索引的列表，包括开始位置，
13     # 几条条目，按字节有多长
14     self.postings_dict[term] = (start_position_in_index_file,
15                                len(postings_list),
16                                length_in_bytes_of_postings_list)

```

### 2.4.2.2 invert\_write()

将解析得到的td\_pairs转换成倒排表，并使用 `InvertedIndexWriter` 类将其写入磁盘。

```

1  def invert_write(self, td_pairs, index):
2      td_dict = collections.defaultdict(list)
3      for t, d in td_pairs:
4          td_dict[t].append(d)
5          # 以键值排序
6          for t in sorted(td_dict.keys()):
7              # 每一个term对应着一个posting_list
8              p_list = sorted(td_dict[t])
9              # 在索引中加入这个tuple
10             index.append(t, sorted(p_list))

```

## 2.4.3 合并

### 2.4.3.1 基础函数

```

1  def _initialization_hook(self):
2      # 初始化指针，一开始在文件的开头
3      self.curr_term_pos = 0
4
5  def __next__(self):
6
7      if self.curr_term_pos >= len(self.terms):
8          raise StopIteration
9

```

```

10     # term
11     term = self.terms[self.curr_term_pos]
12     self.curr_term_pos += 1
13     # 对于每个term，都有一块index的空间。
14     start_position, n_postings, length_in_bytes =
self.postings_dict[term]
15     # offset=start_position
16     self.index_file.seek(start_position)
17     # 读取相应长度的数据，就是postings_list
18     postings_list = self.postings_encoding.decode(
19         self.index_file.read(length_in_bytes))
20     return term, postings_list

```

### 2.4.3.2 merge()

把相同term后面跟着的postings\_list整合到一起。

- indices: 指针列表，每个指针指向其中一个小的索引
- merged\_index: 结果

```

1  def merge(self, indices, merged_index):
2      last_term = last_posting = None
3
4      # 先把indices合并排序，得到一个 term & postings 对的列表（排好序的）。
5      for curr_term, curr_postings in heapq.merge(*indices):
6          # 如果和之前的不相等也就是说之前的那个没有可以合并得了。
7          if curr_term != last_term:
8              # 如果不是第一个条目
9              if last_term:
10                 # 在结果中加上这一term及其条目，接着循环。
11                 last_posting = list(sorted(set(last_posting)))
12                 merged_index.append(last_term, last_posting)
13                 last_term = curr_term
14                 last_posting = curr_postings
15             else:
16                 # 如果这两个需要合并，就把这个条目加进去。
17                 last_posting += curr_postings
18             if last_term:
19                 last_posting = list(sorted(set(last_posting)))
20                 merged_index.append(last_term, last_posting)

```

## 3 布尔联合检索 (10%)

## 3.1 \_get\_postings\_list()

给定一个term，返回对应的条目列表。（不会遍历整个index文件）

```
1 def _get_postings_list(self, term):
2     # 开始地址，posting的条数，字节数
3     start_position, n_postings, length_in_bytes =
4     self.postings_dict[term]
5     self.index_file.seek(start_position)
6     return self.postings_encoding.decode(
7         self.index_file.read(length_in_bytes))
```

## 3.2 sorted\_intersect()

对两个排序好的list做Intersect，最后返回排好序的结果。

```
1 def sorted_intersect(list1, list2):
2     idx1 = idx2 = 0
3     intersect = []
4     while idx1 < len(list1) and idx2 < len(list2):
5         if list1[idx1] < list2[idx2]:
6             idx1 += 1
7         elif list2[idx2] < list1[idx1]:
8             idx2 += 1
9         else:
10            intersect.append(list1[idx1])
11            idx1 += 1
12            idx2 += 1
13    return intersect
```

## 3.3 retrieve()

查找同时符合几个条件的文件

- query: 用空格分隔的token list (string 类型)

返回doc list

```
1 def retrieve(self, query):
2     if len(self.term_id_map) == 0 or len(self.doc_id_map) == 0:
3         self.load()
4
5     with InvertedIndexMapper(self.index_name,
6         directory=self.output_dir,
```

```

6             postings_encoding=
7             self.postings_encoding) as mapper:
8         result = None
9         # 先把query以空格为分界分成list，然后遍历
10        for term in query.split():
11            term_id = self.term_id_map.str_to_id.get(term)
12            if not term_id:
13                return []
14            r = mapper[term_id]
15            if result is None:
16                result = r
17            else:
18                result = sorted_intersect(result, r)
19        return [self.doc_id_map[r] for r in result]

```

## 4 索引压缩 (30%)

```

1  class CompressedPostings:
2      @staticmethod
3      def encode_int(gap):
4          ret = [(gap & 0x7f) | 0x80]
5          gap >>= 7
6          while gap != 0:
7              ret.insert(0, gap & 0x7f)
8              gap >>= 7
9          return ret
10
11     @staticmethod
12     def encode(postings_list):
13         # 初始化结果列表
14         encoded_postings_list = []
15         encoded_postings_list += CompressedPostings.encode_int(
16             postings_list[0])
17         # 对每一个都做压缩
18         for i in range(1, len(postings_list)):
19             encoded_postings_list += CompressedPostings.encode_int(
20                 postings_list[i] - postings_list[i-1])
21         return array.array('B', encoded_postings_list).tobytes()
22
23
24     @staticmethod
25     def decode(encoded_postings_list):
26         decoded_postings_list = array.array('B')
27         decoded_postings_list.frombytes(encoded_postings_list)
28
29         postings_list = []
30         base, n = 0, len(decoded_postings_list)

```

```

31         idx = 0
32         while idx < n:
33             gap = 0
34             while idx < n and (decoded_postings_list[idx] & 0x80) ==
0:
35                 gap = (gap << 7) | (decoded_postings_list[idx] &
0x7f)
36                 idx += 1
37             gap = (gap << 7) | (decoded_postings_list[idx] & 0x7f)
38             idx += 1
39
40             posting = base + gap
41             postings_list.append(posting)
42             base = posting
43         return postings_list

```

## 5 额外的编码方式 (10%)

**gamma-encoding**: 先去掉1, 然后还剩n bit, 然后求出n的一元码, 然后把这俩拼在一起。

```

1  class ECompressedPostings:
2      def encode_int(gap):
3          if gap == 0 or gap == 1:
4              return '0'
5          ret = '1' * int(log(gap, 2)) + '0' + bin(gap)[3:]
6          print(ret)
7          return ret
8
9      @staticmethod
10     def encode(postings_list):
11         encoded_postings_list = ''
12         encoded_postings_list += ECompressedPostings.encode_int(
13             postings_list[0] - (-1))
14         for i in range(1, len(postings_list)):
15             encoded_postings_list +=
ECompressedPostings.encode_int(
16                 postings_list[i] - postings_list[i - 1])
17         print(encoded_postings_list)
18         return array.array('B', [int(encoded_postings_list[x:x+8],
2) for x in range(0, len(encoded_postings_list), 8)]).tobytes()
19
20
21     @staticmethod
22     def decode(encoded_postings_list):
23         decoded_bytes_list = array.array('B')
24         decoded_bytes_list.frombytes(encoded_postings_list)

```

```

25
26     decoded_postings_list = ''.join([bin(x)[2:].zfill(8)
27                                     for x in
decoded_bytes_list])
28     decoded_postings_list = decoded_postings_list[:-7] +
bin(decoded_bytes_list[-1])[2:]
29     print(decoded_postings_list)
30
31     postings_list = []
32     base, idx, n = -1, 0, len(decoded_postings_list)
33     while idx < n:
34         length = 0
35         while idx < n and decoded_postings_list[idx] == '1':
36             length += 1
37             idx += 1
38         if idx < n:
39             # '111 ... 1(length)0xxx ... x(length)', length maybe 0
40             idx = idx + 1 + length
41             gap = int('1' + decoded_postings_list[idx-length :
idx], 2)
42             print(idx, gap)
43             posting = base + gap
44             postings_list.append(posting)
45             base = posting
46     return postings_list

```