

自然语言处理

第三次编程作业

西安交通大学, 数学与统计学院, 强基数学 002

马煜璇^a, 吴天阳^b

2204220461^a, 2204210460^b

2022 年 12 月 3 日

目录

1	实验目的	3
2	实验原理	3
2.1	正向最大匹配算法	3
2.2	逆向最大匹配算法	3
2.3	双向最大匹配算法	4
3	实验步骤与结果分析	4
3.1	读入词库	4
3.2	正向最大匹配算法	5
3.3	逆向最大匹配算法	5
3.4	双向最大匹配算法	6
4	结论与讨论	6

1 实验目的

利用最大匹配算法对汉语进行分词操作，并对结果进行分析。我们分别实现下面三种算法：

1. 正向最大匹配算法 (FMM)
2. 反向最大匹配算法 (BMM)
3. 双向最大匹配算法

2 实验原理

2.1 正向最大匹配算法

正向最大匹配算法 (FMM) 匹配的基本思想是：假设词典中最大词条所含的汉字个数为 n 个，取待处理字符串的前 n 个字作为匹配字段，查找分词词典。若词典中含有该词，则匹配成功，分出该词，然后从被比较字符串的 $n+1$ 处开始再取 n 个字组成的字段重新在词典中匹配；如果没有匹配成功，则将这 n 个字组成的字段的最后一位剔除，用剩下的 $n-1$ 个字组成的字段在词典中进行匹配，如此进行下去，直到切分成功为止。

例如，待处理字符串为“汉字多为表意文字”，取字符串“汉语多为表”与词典进行比较，没有与之对应的词，去除“表”字，用字段“汉语多为”进行匹配，直至匹配到“汉语”为止，再取字符串“多为表意”，循环到切分出“文字”一词。

目前，正向最大匹配方法作为一种基本的方法已被肯定下来，但是由于错误比较大，一般不单独使用。如字符串“处理机器发生的故障”，在正向最大匹配方法中会出现歧义切分，该字符串被分为：处理机、发生、故障，但是使用逆向匹配就能得到有效的切分。

2.2 逆向最大匹配算法

逆向最大匹配算法 (RMM) 的分词原理和过程与正向最大匹配相似，区别在于前者从文章或者句子 (字串) 的末尾开始切分，若不成功则减去最前面的一个字。

比如对于字符串“处理机器发生的故障”，第一步，从字串的右边取长度以步长为单位的字段“发生的故障”在词典中进行匹配，匹配不成功，再取字段“生的故障”进行匹配，依次匹配，直到分出“故障”一词，最终使用 RMM 方法切分的结果为：故障、发生、机器、处理。该方法要求配备逆序词典。

一般来说根据汉语词汇构成的特点，从理论上说明了逆向匹配的精确度高于正向匹配，汉语语句的特点一般中心语偏后。有研究数据，单纯使用正向最大匹配的错误率为 $1/169$ ，单纯使用逆向最大匹配的错误率为 $1/245$ 。

2.3 双向最大匹配算法

FMM 和 BMM 两种算法都分词一遍，然后根据大颗粒度词越多越好，非词典词和单字词越少越好的原则，选取其中一种分词结果输出。

3 实验步骤与结果分析

3.1 读入词库

我们从https://github.com/fxsjy/jieba/blob/master/extra_dict/dict.txt.big下载了词库，并读取词库，词库大小为 583280。

```
1 # 1.词库读取
2 dic = set()
3 with open('dict.txt_2.big', encoding='utf-8') as file:
4     for s in file.readlines():
5         word = s.split()[0]
6         dic.add(word)
7 print('词库大小', len(dic))
```

同时，我们查找了词库中最长的词语与最长词的长度，长度为 16。

```
1 #词库中的最长词
2 counter = collections.Counter()
3 for w in dic:
4     counter[w] = len(w)
5 print(counter.most_common()[:10])
```

词库中最长的词语与长度为

```
1 [('第九届全国人民代表大会常务委员会', 16), ('侵華日軍南京大屠殺遇難同胞紀念館', 16),
   ('侵华日军南京大屠杀遇难同胞纪念馆', 16), ('第九屆全國人民代表大會常務委員會',
   16), ('劳动和社会保障部职业技能鉴定中心', 16),
   ('八千一百三十七萬七千二百三十六口', 16), ('外交部駐澳門特別行政區特派員公署',
   16), ('外交部駐香港特別行政區特派員公署', 16),
   ('外交部驻香港特别行政区特派员公署', 16), ('八千一百三十七万七千二百三十六口',
   16)]
```

计算文本长度

```
1 #原文本长度
```

```
2 length = len(text) # 获取文本长度
3 print('原文本长度', length, '\n')
```

判断是否为词语

```
1 def check(word): # 判断是否为词语
2     return word in dic
```

3.2 正向最大匹配算法

正向最大匹配算法从正向取词，即从左往右取词，取词最大长度为词典中最长词的长度16，每次右边减一个字，直到词典中存在或剩下1个单字。

```
1 pos_div= []
2 pos_single = 0
3 i = 0 # 左指针
4 while i < length:
5     for j in range(min(i+16, length), i, -1): # 右指针
6         if check(text[i:j]):
7             pos_div.append(text[i:j])
8             break
9     if i + 1 == j: # 记录单个字词个数
10        pos_single += 1
11    i = j
12 print('正向最大匹配法长度: ', len(pos_div))
13 print('单个字个数: ', pos_single)
```

3.3 逆向最大匹配算法

反向即从右往左取词，其他逻辑和正向相同。

```
1 neg_div= []
2 neg_single = 0
3 j = length # 右指针
4 while j > 0:
5     for i in range(max(j-16, 0), j): # 左指针
6         if check(text[i:j]):
7             neg_div.append(text[i:j])
8             break
```

```
9     if i + 1 == j: # 记录单个字词个数
10         neg_single += 1
11     j = i
12 neg_div = neg_div[::-1]
13 print('逆向最大匹配法长度: ', len(neg_div))
14 print('单个字个数: ', neg_single, '\n')
```

3.4 双向最大匹配算法

FMM 和 BMM 两种算法都分词一遍，选取其中一种分词结果输出。选择标准：

1. 首先看两种方法结果的分词数，分词数越少越好；
2. 分词数相同的情况下，看单个词的数量，越少越好。

```
1 best_div = None
2 info = [len(pos_div)-len(neg_div), pos_single-neg_single]
3 # 优先选取分词个数少的，再选取单个词少的
4 if info[0] < 0 or (info[0] == 0 and info[1] < 0):
5     best_div = pos_div
6     print('选择正向')
7 elif info[0] == info[1] == 0: # 两个数目都相同
8     best_div = pos_div
9     print('两者相同')
10 else:
11     print('选择逆向')
12     best_div = neg_div
13 print(best_div)
```

4 结论与讨论

我们对朱自清先生的文章《荷塘月色》片段进行分词，文本长度为 787，分别利用正向最大匹配算法和逆向最大匹配算法对其进行分词，得到分词结果如下：

```
1 正向最大匹配法长度:  448
2 单个字个数:  355
3 逆向最大匹配法长度:  447
4 单个字个数:  354
```

根据双向最大匹配算法的选择标准，选择逆向，原文分词结果为

[‘路上’，‘只’，‘我’，‘一’，‘个人’，‘背着手’，‘踱’，‘着’，‘这’，‘一’，‘片’，‘天’，‘地’，‘好’，‘像’
是’，‘我’，‘的’，‘我’，‘也’，‘像’，‘超出’，‘了’，‘平常’，‘的’，‘自己’，‘到’，‘了’，‘另’，‘一’，
世界’，‘里’，‘我’，‘爱’，‘热闹’，‘也’，‘爱’，‘冷静’，‘爱’，‘群居’，‘也’，‘爱’，‘独处’，‘像’，
今’，‘晚上’，‘一’，‘个人’，‘在’，‘这’，‘苍’，‘茫’，‘的’，‘月’，‘下’，‘什’，‘么’，‘都’，‘可’，‘以’，‘想’，
什’，‘么’，‘都’，‘可’，‘以’，‘不’，‘想’，‘便’，‘觉’，‘是’，‘个’，‘自’，‘由’，‘的’，‘人’，‘白’，‘天’，‘里’，‘一’，
定要’，‘做’，‘的’，‘事’，‘一’，‘定’，‘要’，‘说’，‘的’，‘话’，‘现’，‘在’，‘都’，‘可’，‘不’，‘理’，‘这’，‘是’，‘独’，
处’，‘的’，‘妙’，‘处’，‘我’，‘且’，‘受’，‘用’，‘这’，‘无’，‘边’，‘的’，‘荷’，‘香’，‘月’，‘色’，‘好’，‘了’，‘曲’，‘曲’，
折’，‘折’，‘的’，‘荷’，‘塘’，‘上’，‘面’，‘弥’，‘望’，‘的’，‘是’，‘田’，‘田’，‘的’，‘叶’，‘子’，‘叶’，‘子’，‘出’，‘水’，‘很’，
高’，‘像’，‘亭’，‘亭’，‘的’，‘舞’，‘女’，‘的’，‘裙’，‘层’，‘层’，‘的’，‘叶’，‘子’，‘中’，‘间’，‘零’，‘星’，‘地’，‘点’，
缀’，‘着’，‘些’，‘白’，‘花’，‘有’，‘袅’，‘娜’，‘地’，‘开’，‘着’，‘的’，‘有’，‘羞’，‘涩’，‘地’，‘打’，‘着’，‘朵’，‘儿’，
的’，‘正’，‘如’，‘一’，‘粒’，‘粒’，‘的’，‘明’，‘珠’，‘又’，‘如’，‘碧’，‘天’，‘里’，‘的’，‘星’，‘星’，‘又’，‘如’，
‘刚’，‘出’，‘浴’，‘的’，‘美’，‘人’，‘微’，‘风’，‘过’，‘处’，‘送’，‘来’，‘缕’，‘缕’，‘清’，‘香’，‘仿’，‘佛’，‘远’，‘处’，
高’，‘楼’，‘上’，‘渺’，‘茫’，‘的’，‘歌’，‘声’，‘似’，‘的’，‘这’，‘时’，‘候’，‘叶’，‘子’，‘与’，‘花’，‘也’，‘有’，‘一’，‘丝’，
的’，‘颤’，‘动’，‘像’，‘闪’，‘电’，‘般’，‘霎’，‘时’，‘传’，‘过’，‘荷’，‘塘’，‘的’，‘那’，‘边’，‘去’，‘了’，‘叶’，‘子’，‘本’，
是’，‘肩’，‘并’，‘肩’，‘密’，‘密’，‘地’，‘挨’，‘着’，‘这’，‘便’，‘宛’，‘然’，‘有’，‘了’，‘一’，‘道’，‘凝’，‘碧’，‘的’，
波’，‘痕’，‘叶’，‘子’，‘底’，‘下’，‘是’，‘脉’，‘脉’，‘的’，‘流’，‘水’，‘遮’，‘住’，‘了’，‘不’，‘能’，‘见’，‘一’，‘些’，
颜’，‘色’，‘而’，‘叶’，‘子’，‘却’，‘更’，‘见’，‘风’，‘致’，‘了’，‘月’，‘光’，‘如’，‘流’，‘水’，‘一’，‘般’，‘静’，‘静’，
地’，‘泻’，‘在’，‘这’，‘一’，‘片’，‘叶’，‘子’，‘和’，‘花’，‘上’，‘薄’，‘薄’，‘的’，‘青’，‘雾’，‘浮’，‘起’，‘在’，‘荷’，‘塘’，‘里’，
叶’，‘子’，‘和’，‘花’，‘仿’，‘佛’，‘在’，‘牛’，‘乳’，‘中’，‘洗’，‘过’，‘一’，‘样’，‘又’，‘像’，‘笼’，‘着’，‘轻’，‘纱’，
的’，‘梦’，‘虽’，‘然’，‘是’，‘满’，‘月’，‘天’，‘上’，‘却’，‘有’，‘一’，‘层’，‘淡’，‘淡’，‘的’，‘云’，‘所’，‘以’，‘不’，‘能’，
朗’，‘照’，‘但’，‘我’，‘以’，‘为’，‘这’，‘恰’，‘是’，‘到’，‘了’，‘好’，‘处’，‘酣’，‘眠’，‘固’，‘不’，‘可’，‘少’，
‘小’，‘睡’，‘也’，‘别’，‘有’，‘风’，‘味’，‘的’，‘月’，‘光’，‘是’，‘隔’，‘了’，‘树’，‘照’，‘过’，‘来’，‘的’，‘高’，‘处’，‘从’，
生’，‘的’，‘灌’，‘木’，‘落’，‘下’，‘参’，‘差’，‘的’，‘斑’，‘驳’，‘的’，‘黑’，‘影’，‘峭’，‘楞’，‘楞’，‘如’，‘鬼’，‘一’，
般’，‘弯’，‘弯’，‘的’，‘杨’，‘柳’，‘的’，‘稀’，‘疏’，‘的’，‘倩’，‘影’，‘却’，‘又’，‘像’，‘是’，‘画’，‘在’，‘荷’，‘叶’，
上’，‘塘’，‘中’，‘的’，‘月’，‘色’，‘并’，‘不’，‘均’，‘匀’，‘但’，‘光’，‘与’，‘影’，‘有’，‘着’，‘和’，‘谐’，‘的’，‘旋’，
律’，‘如’，‘梵’，‘婀’，‘玲’，‘英’，‘语’，‘小’，‘提’，‘琴’，‘的’，‘译’，‘音’，‘上’，‘奏’，‘着’，‘的’，‘名’，‘曲’，‘荷’，
塘’，‘的’，‘四’，‘面’，‘远’，‘远’，‘近’，‘近’，‘高’，‘高’，‘低’，‘低’，‘都’，‘是’，‘树’，‘而’，‘杨’，‘柳’，‘最’，‘多’，‘这’，‘些’，
树’，‘将’，‘一’，‘片’，‘荷’，‘塘’，‘重’，‘重’，‘围’，‘住’，‘只’，‘在’，‘小’，‘路’，‘一’，‘旁’，‘漏’，‘着’，‘几’，‘段’，‘空’，‘隙’，
‘像’，‘是’，‘特’，‘为’，‘月’，‘光’，‘留’，‘下’，‘的’，‘树’，‘色’，‘一’，‘例’，‘是’，‘阴’，‘阴’，‘的’，‘乍’，‘看’，‘像’，‘一’，
团’，‘烟’，‘雾’，‘但’，‘杨’，‘柳’，‘的’，‘丰’，‘姿’，‘便’，‘在’，‘烟’，‘雾’，‘里’，‘也’，‘辨’，‘得’，‘出’，‘树’，‘梢’，
上’，‘隐’，‘隐’，‘约’，‘约’，‘的’，‘是’，‘一’，‘带’，‘远’，‘山’，‘只’，‘有’，‘些’，‘大’，‘意’，‘罢’，‘了’，‘树’，‘缝’，‘里’，‘也’，
‘漏’，‘着’，‘一’，‘两’，‘点’，‘路’，‘灯’，‘光’，‘没’，‘精’，‘打’，‘采’，‘的’，‘是’，‘渴’，‘睡’，‘人’，‘的’，‘眼’，‘这’，‘时’，
候’，‘最’，‘热’，‘闹’，‘的’，‘要’，‘数’，‘树’，‘上’，‘的’，‘蝉’，‘声’，‘与’，‘水’，‘里’，‘的’，‘蛙’，‘声’，‘但’，‘热’，
闹’，‘是’，‘它’，‘们’，‘的’，‘我’，‘什’，‘么’，‘也’，‘没’，‘有’]

附录

```
1 import collections
2
3 def solve(text):
4     # 1.词库读取
5     dic = set()
6     with open('dict.txt_2.big', encoding='utf-8') as file:
7         for s in file.readlines():
8             word = s.split()[0]
9             dic.add(word)
10    print('词库大小', len(dic))
11
12    #词库中的最长词
13    counter = collections.Counter()
14    for w in dic:
15        counter[w] = len(w)
16    print(counter.most_common()[:10])
17
18    #原文本长度
19    length = len(text) # 获取文本长度
20    print('原文本长度', length, '\n')
21
22    def check(word): # 判断是否为词语
23        return word in dic
24
25    # 2.正向最大匹配法
26    pos_div= []
27    pos_single = 0
28    i = 0 # 左指针
29    while i < length:
30        for j in range(min(i+16, length), i, -1): # 右指针
31            if check(text[i:j]):
32                pos_div.append(text[i:j])
33                break
34            if i + 1 == j: # 记录单个字词个数
35                pos_single += 1
36            i = j
```



```

37 print('正向最大匹配法长度: ', len(pos_div))
38 print('单个字个数: ', pos_single)
39
40 # 3.逆向最大匹配法
41 neg_div= []
42 neg_single = 0
43 j = length # 右指针
44 while j > 0:
45     for i in range(max(j-16, 0), j): # 左指针
46         if check(text[i:j]):
47             neg_div.append(text[i:j])
48             break
49         if i + 1 == j: # 记录单个字词个数
50             neg_single += 1
51         j = i
52     neg_div = neg_div[::-1]
53     print('逆向最大匹配法长度: ', len(neg_div))
54     print('单个字个数: ', neg_single, '\n')
55
56 # 4.双向最大匹配法
57 best_div = None
58 info = [len(pos_div)-len(neg_div), pos_single-neg_single]
59 # 优先选取分词个数少的, 再选取单个词少的
60 if info[0] < 0 or (info[0] == 0 and info[1] < 0):
61     best_div = pos_div
62     print('选择正向')
63 elif info[0] == info[1] == 0: # 两个数目都相同
64     best_div = pos_div
65     print('两者相同')
66 else:
67     print('选择逆向')
68     best_div = neg_div
69     print(best_div)
70
71 if __name__ == '__main__':
72     # with open('text.txt', encoding='utf-8') as file:
73     #     long_text = file.read()
74     text = '路上只我一个人, 背着手踱着。这一片天地好像是我的...' # 设定目标文本

```

